



# NVML

vR550 | March 2024

## Reference Manual



# TABLE OF CONTENTS

Chapter 1. NVML API Reference.....	1
Chapter 2. Known Issues.....	3
Chapter 3. Change Log.....	4
Chapter 4. Modules.....	17
4.1. Device Structs.....	18
nvmlPciInfoExt_v1_t.....	19
nvmlPciInfo_t.....	19
nvmlEccErrorCounts_t.....	19
nvmlUtilization_t.....	19
nvmlMemory_t.....	19
nvmlMemory_v2_t.....	19
nvmlBAR1Memory_t.....	19
nvmlProcessInfo_v1_t.....	19
nvmlProcessInfo_t.....	19
nvmlProcessDetail_v1_t.....	19
nvmlProcessDetailList_v1_t.....	19
nvmlC2cModelInfo_v1_t.....	19
nvmlRowRemapperHistogramValues_t.....	19
nvmlNvLinkUtilizationControl_t.....	19
nvmlBridgeChipInfo_t.....	19
nvmlBridgeChipHierarchy_t.....	19
nvmlValue_t.....	19
nvmlSample_t.....	19
nvmlViolationTime_t.....	20
nvmlBridgeChipType_t.....	20
nvmlNvLinkUtilizationCountUnits_t.....	20
nvmlNvLinkUtilizationCountPktTypes_t.....	20
nvmlNvLinkCapability_t.....	20
nvmlNvLinkErrorCounter_t.....	21
nvmlIntNvLinkDeviceType_t.....	21
nvmlGpuTopologyLevel_t.....	21
nvmlSamplingType_t.....	22
nvmlPcieUtilCounter_t.....	22
nvmlValueType_t.....	22
nvmlPerfPolicyType_t.....	23
NVML_VALUE_NOT_AVAILABLE.....	23
NVML_DEVICE_PCI_BUS_ID_BUFFER_SIZE.....	23
NVML_DEVICE_PCI_BUS_ID_BUFFER_V2_SIZE.....	24
NVML_DEVICE_PCI_BUS_ID_LEGACY_FMT.....	24
NVML_DEVICE_PCI_BUS_ID_FMT.....	24

NVML_DEVICE_PCI_BUS_ID_FMT_ARGS.....	24
nvmlProcessDetailList_v1.....	24
NVML_NVLINK_MAX_LINKS.....	24
NVML_MAX_PHYSICAL_BRIDGE.....	24
4.2. Device Enums.....	24
nvmlClkMonFaultInfo_t.....	25
nvmlClkMonStatus_t.....	25
nvmlEnableState_t.....	25
nvmlBrandType_t.....	25
nvmlTemperatureThresholds_t.....	26
nvmlTemperatureSensors_t.....	26
nvmlComputeMode_t.....	26
nvmlMemoryErrorType_t.....	27
nvmlEccCounterType_t.....	27
nvmlClockType_t.....	27
nvmlClockId_t.....	28
nvmlDriverModel_t.....	28
nvmlPstates_t.....	28
nvmlGpuOperationMode_t.....	29
nvmlInforomObject_t.....	30
nvmlReturn_t.....	30
nvmlMemoryLocation_t.....	32
nvmlPageRetirementCause_t.....	32
nvmlRestrictedAPI_t.....	33
nvmlFlagDefault.....	33
nvmlFlagForce.....	33
MAX_CLK_DOMAINS.....	33
nvmlEccBitType_t.....	33
NVML_SINGLE_BIT_ECC.....	33
NVML_DOUBLE_BIT_ECC.....	33
4.3. Field Value Enums.....	34
nvmlFieldValue_t.....	34
NVML_FI_DEV_ECC_CURRENT.....	34
NVML_FI_DEV_ECC_PENDING.....	34
NVML_FI_DEV_ECC_SBE_VOL_TOTAL.....	34
NVML_FI_DEV_ECC_DBE_VOL_TOTAL.....	34
NVML_FI_DEV_ECC_SBE_AGG_TOTAL.....	34
NVML_FI_DEV_ECC_DBE_AGG_TOTAL.....	34
NVML_FI_DEV_ECC_SBE_VOL_L1.....	34
NVML_FI_DEV_ECC_DBE_VOL_L1.....	34
NVML_FI_DEV_ECC_SBE_VOL_L2.....	35
NVML_FI_DEV_ECC_DBE_VOL_L2.....	35
NVML_FI_DEV_ECC_SBE_VOL_DEV.....	35

NVML_FI_DEV_ECC_DBE_VOL_DEV.....	35
NVML_FI_DEV_ECC_SBE_VOL_REG.....	35
NVML_FI_DEV_ECC_DBE_VOL_REG.....	35
NVML_FI_DEV_ECC_SBE_VOL_TEX.....	35
NVML_FI_DEV_ECC_DBE_VOL_TEX.....	35
NVML_FI_DEV_ECC_DBE_VOL_CBU.....	35
NVML_FI_DEV_ECC_SBE_AGG_L1.....	35
NVML_FI_DEV_ECC_DBE_AGG_L1.....	35
NVML_FI_DEV_ECC_SBE_AGG_L2.....	36
NVML_FI_DEV_ECC_DBE_AGG_L2.....	36
NVML_FI_DEV_ECC_SBE_AGG_DEV.....	36
NVML_FI_DEV_ECC_DBE_AGG_DEV.....	36
NVML_FI_DEV_ECC_SBE_AGG_REG.....	36
NVML_FI_DEV_ECC_DBE_AGG_REG.....	36
NVML_FI_DEV_ECC_SBE_AGG_TEX.....	36
NVML_FI_DEV_ECC_DBE_AGG_TEX.....	36
NVML_FI_DEV_ECC_DBE_AGG_CBU.....	36
NVML_FI_DEV_RETIRED_SBE.....	36
NVML_FI_DEV_RETIRED_DBE.....	36
NVML_FI_DEV_RETIRED_PENDING.....	37
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L0.....	37
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L1.....	37
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L2.....	37
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L3.....	37
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L4.....	37
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L5.....	37
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_TOTAL.....	37
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L0.....	38
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L1.....	38
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L2.....	38
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L3.....	38
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L4.....	38
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L5.....	38
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_TOTAL.....	38
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L0.....	38
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L1.....	39
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L2.....	39
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L3.....	39
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L4.....	39
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L5.....	39
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_TOTAL.....	39
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L0.....	39
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L1.....	39

NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L2.....	40
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L3.....	40
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L4.....	40
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L5.....	40
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_TOTAL.....	40
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L0.....	40
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L1.....	40
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L2.....	40
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L3.....	40
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L4.....	41
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L5.....	41
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_TOTAL.....	41
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L0.....	41
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L1.....	41
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L2.....	41
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L3.....	41
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L4.....	41
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L5.....	41
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_TOTAL.....	41
NVML_FI_DEV_PERF_POLICY_POWER.....	41
NVML_FI_DEV_PERF_POLICY_THERMAL.....	42
NVML_FI_DEV_PERF_POLICY_SYNC_BOOST.....	42
NVML_FI_DEV_PERF_POLICY_BOARD_LIMIT.....	42
NVML_FI_DEV_PERF_POLICY_LOW_UTILIZATION.....	42
NVML_FI_DEV_PERF_POLICY_RELIABILITY.....	42
NVML_FI_DEV_PERF_POLICY_TOTAL_APP_CLOCKS.....	42
NVML_FI_DEV_PERF_POLICY_TOTAL_BASE_CLOCKS.....	42
NVML_FI_DEV_MEMORY_TEMP.....	42
NVML_FI_DEV_TOTAL_ENERGY_CONSUMPTION.....	42
NVML_FI_DEV_NVLINK_SPEED_MBPS_L0.....	42
NVML_FI_DEV_NVLINK_SPEED_MBPS_L1.....	43
NVML_FI_DEV_NVLINK_SPEED_MBPS_L2.....	43
NVML_FI_DEV_NVLINK_SPEED_MBPS_L3.....	43
NVML_FI_DEV_NVLINK_SPEED_MBPS_L4.....	43
NVML_FI_DEV_NVLINK_SPEED_MBPS_L5.....	43
NVML_FI_DEV_NVLINK_SPEED_MBPS_COMMON.....	43
NVML_FI_DEV_NVLINK_LINK_COUNT.....	43
NVML_FI_DEV_RETIRED_PENDING_SBE.....	43
NVML_FI_DEV_RETIRED_PENDING_DBE.....	43
NVML_FI_DEV_PCIE_REPLAY_COUNTER.....	43
NVML_FI_DEV_PCIE_REPLAY_ROLLOVER_COUNTER.....	43
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L6.....	44
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L7.....	44

NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L8.....	44
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L9.....	44
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L10.....	44
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L11.....	44
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L6.....	44
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L7.....	44
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L8.....	45
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L9.....	45
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L10.....	45
NVML_FI_DEV_NVLINK_CRC_DATA_ERROR_COUNT_L11.....	45
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L6.....	45
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L7.....	45
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L8.....	45
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L9.....	45
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L10.....	46
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L11.....	46
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L6.....	46
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L7.....	46
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L8.....	46
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L9.....	46
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L10.....	46
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L11.....	47
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L6.....	47
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L7.....	47
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L8.....	47
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L9.....	47
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L10.....	47
NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L11.....	47
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L6.....	47
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L7.....	47
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L8.....	47
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L9.....	47
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L10.....	48
NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L11.....	48
NVML_FI_DEV_NVLINK_SPEED_MBPS_L6.....	48
NVML_FI_DEV_NVLINK_SPEED_MBPS_L7.....	48
NVML_FI_DEV_NVLINK_SPEED_MBPS_L8.....	48
NVML_FI_DEV_NVLINK_SPEED_MBPS_L9.....	48
NVML_FI_DEV_NVLINK_SPEED_MBPS_L10.....	48
NVML_FI_DEV_NVLINK_SPEED_MBPS_L11.....	48
NVML_FI_DEV_NVLINK_THROUGHPUT_DATA_TX.....	48
NVML_FI_DEV_NVLINK_THROUGHPUT_DATA_RX.....	49
NVML_FI_DEV_NVLINK_THROUGHPUT_RAW_TX.....	49

NVML_FI_DEV_NVLINK_THROUGHPUT_RAW_RX.....	49
NVML_FI_DEV_REMAPPED_COR.....	49
NVML_FI_DEV_REMAPPED_UNC.....	49
NVML_FI_DEV_REMAPPED_PENDING.....	49
NVML_FI_DEV_REMAPPED_FAILURE.....	49
NVML_FI_DEV_NVLINK_REMOTE_NVLINK_ID.....	49
NVML_FI_DEV_NVSWITCH_CONNECTED_LINK_COUNT.....	49
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L0.....	50
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L1.....	50
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L2.....	50
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L3.....	50
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L4.....	50
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L5.....	50
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L6.....	50
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L7.....	50
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L8.....	51
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L9.....	51
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L10.....	51
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_L11.....	51
NVML_FI_DEV_NVLINK_ECC_DATA_ERROR_COUNT_TOTAL.....	51
NVML_FI_DEV_NVLINK_ERROR_DL_REPLAY.....	51
NVML_FI_DEV_NVLINK_ERROR_DL_RECOVERY.....	51
NVML_FI_DEV_NVLINK_ERROR_DL_CRC.....	51
NVML_FI_DEV_NVLINK_GET_SPEED.....	52
NVML_FI_DEV_NVLINK_GET_STATE.....	52
NVML_FI_DEV_NVLINK_GET_VERSION.....	52
NVML_FI_DEV_NVLINK_GET_POWER_STATE.....	52
NVML_FI_DEV_NVLINK_GET_POWER_THRESHOLD.....	52
NVML_FI_DEV_PCIE_L0_TO_RECOVERY_COUNTER.....	52
NVML_FI_DEV_C2C_LINK_COUNT.....	52
NVML_FI_DEV_C2C_LINK_GET_STATUS.....	52
NVML_FI_DEV_C2C_LINK_GET_MAX_BW.....	52
NVML_FI_DEV_POWER_AVERAGE.....	52
NVML_FI_DEV_POWER_INSTANT.....	53
NVML_FI_DEV_POWER_MIN_LIMIT.....	53
NVML_FI_DEV_POWER_MAX_LIMIT.....	53
NVML_FI_DEV_POWER_DEFAULT_LIMIT.....	53
NVML_FI_DEV_POWER_CURRENT_LIMIT.....	53
NVML_FI_DEV_ENERGY.....	53
NVML_FI_DEV_POWER_REQUESTED_LIMIT.....	53
NVML_FI_DEV_TEMPERATURE_SHUTDOWN_TLIMIT.....	53
NVML_FI_DEV_TEMPERATURE_SLOWDOWN_TLIMIT.....	54
NVML_FI_DEV_TEMPERATURE_MEM_MAX_TLIMIT.....	54

NVML_FI_DEV_TEMPERATURE_GPU_MAX_TLIMIT.....	54
NVML_FI_DEV_IS_MIG_MODE_INDEPENDENT_MIG_QUERY_CAPABLE.....	54
NVML_FI_MAX.....	54
4.4. Unit Structs.....	54
nvmlHwbcEntry_t.....	55
nvmlLedState_t.....	55
nvmlUnitInfo_t.....	55
nvmlPSUInfo_t.....	55
nvmlUnitFanInfo_t.....	55
nvmlUnitFanSpeeds_t.....	55
nvmlFanState_t.....	55
nvmlLedColor_t.....	55
4.5. Accounting Statistics.....	55
nvmlAccountingStats_t.....	56
nvmlDeviceGetAccountingMode.....	56
nvmlDeviceGetAccountingStats.....	56
nvmlDeviceGetAccountingPids.....	57
nvmlDeviceGetAccountingBufferSize.....	58
nvmlDeviceSetAccountingMode.....	59
nvmlDeviceClearAccountingPids.....	60
4.6. Encoder Structs.....	61
nvmlEncoderSessionInfo_t.....	61
nvmlEncoderType_t.....	61
4.7. Frame Buffer Capture Structures.....	61
nvmlFBCStats_t.....	61
nvmlFBCSessionInfo_t.....	61
nvmlFBCSessionType_t.....	61
NVML_NVFBC_SESSION_FLAG_DIFFMAP_ENABLED.....	62
NVML_NVFBC_SESSION_FLAG_CLASSIFICATIONMAP_ENABLED.....	62
NVML_NVFBC_SESSION_FLAG_CAPTURE_WITH_WAIT_NO_WAIT.....	62
NVML_NVFBC_SESSION_FLAG_CAPTURE_WITH_WAIT_INFINITE.....	62
NVML_NVFBC_SESSION_FLAG_CAPTURE_WITH_WAIT_TIMEOUT.....	62
4.8. definitions related to the drain state.....	62
nvmlDetachGpuState_t.....	63
nvmlPcieLinkState_t.....	63
4.9. /nvmlDevice definitions related to Confidential Computing.....	63
nvmlSystemConfComputeSettings_v1_t.....	63
nvmlConfComputeMemSizeInfo_t.....	63
NVML_CC_SYSTEM_CPU_CAPS_NONE.....	63
NVML_CC_SYSTEM_GPUS_CC_NOT_CAPABLE.....	63
NVML_CC_SYSTEM_DEVTOOLS_MODE_OFF.....	63
NVML_CC_SYSTEM_ENVIRONMENT_UNAVAILABLE.....	63
NVML_CC_SYSTEM_FEATURE_DISABLED.....	64



NVML_CC_SYSTEM_MULTIGPU_NONE.....	64
NVML_CC_ACCEPTING_CLIENT_REQUESTS_FALSE.....	64
NVML_GPU_CERT_CHAIN_SIZE.....	64
NVML_CC_GPU_CEC_NONCE_SIZE.....	64
4.10. Initialization and Cleanup.....	64
nvmlInit_v2.....	64
nvmlInitWithFlags.....	65
nvmlShutdown.....	66
NVML_INIT_FLAG_NO_GPUS.....	66
NVML_INIT_FLAG_NO_ATTACH.....	66
4.11. Error reporting.....	66
nvmlErrorString.....	66
4.12. Constants.....	67
NVML_DEVICE_INFOROM_VERSION_BUFFER_SIZE.....	67
NVML_DEVICE_UUID_BUFFER_SIZE.....	67
NVML_DEVICE_UUID_V2_BUFFER_SIZE.....	67
NVML_DEVICE_PART_NUMBER_BUFFER_SIZE.....	67
NVML_SYSTEM_DRIVER_VERSION_BUFFER_SIZE.....	67
NVML_SYSTEM_NVML_VERSION_BUFFER_SIZE.....	67
NVML_DEVICE_NAME_BUFFER_SIZE.....	67
NVML_DEVICE_NAME_V2_BUFFER_SIZE.....	67
NVML_DEVICE_SERIAL_BUFFER_SIZE.....	68
NVML_DEVICE_VBIOS_VERSION_BUFFER_SIZE.....	68
4.13. System Queries.....	68
nvmlSystemGetDriverVersion.....	68
nvmlSystemGetNVMLVersion.....	69
nvmlSystemGetCudaDriverVersion.....	69
nvmlSystemGetCudaDriverVersion_v2.....	70
nvmlSystemGetProcessName.....	70
nvmlSystemGetHicVersion.....	71
nvmlSystemGetTopologyGpuSet.....	72
NVML_CUDA_DRIVER_VERSION_MAJOR.....	72
4.14. Unit Queries.....	72
nvmlUnitGetCount.....	73
nvmlUnitGetHandleByIndex.....	73
nvmlUnitGetUnitInfo.....	74
nvmlUnitGetLedState.....	74
nvmlUnitGetPsuInfo.....	75
nvmlUnitGetTemperature.....	76
nvmlUnitGetFanSpeedInfo.....	76
nvmlUnitGetDevices.....	77
4.15. Device Queries.....	78
CPU and Memory Affinity.....	78

nvmlDeviceGetCount_v2.....	78
nvmlDeviceGetAttributes_v2.....	79
nvmlDeviceGetHandleByIndex_v2.....	79
nvmlDeviceGetHandleBySerial.....	81
nvmlDeviceGetHandleByUUID.....	82
nvmlDeviceGetHandleByPciBusId_v2.....	83
nvmlDeviceGetName.....	84
nvmlDeviceGetBrand.....	84
nvmlDeviceGetIndex.....	85
nvmlDeviceGetSerial.....	86
nvmlDeviceGetC2cModelInfoV.....	87
nvmlDeviceGetTopologyCommonAncestor.....	87
nvmlDeviceGetTopologyNearestGpus.....	88
nvmlDeviceGetP2PStatus.....	89
nvmlDeviceGetUUID.....	89
nvmlDeviceGetMinorNumber.....	90
nvmlDeviceGetBoardPartNumber.....	91
nvmlDeviceGetInforomVersion.....	92
nvmlDeviceGetInforomImageVersion.....	93
nvmlDeviceGetInforomConfigurationChecksum.....	94
nvmlDeviceValidateInforom.....	94
nvmlDeviceGetLastBBXFlushTime.....	95
nvmlDeviceGetDisplayMode.....	96
nvmlDeviceGetDisplayActive.....	96
nvmlDeviceGetPersistenceMode.....	97
nvmlDeviceGetPciInfoExt.....	98
nvmlDeviceGetPciInfo_v3.....	98
nvmlDeviceGetMaxPcieLinkGeneration.....	99
nvmlDeviceGetGpuMaxPcieLinkGeneration.....	100
nvmlDeviceGetMaxPcieLinkWidth.....	100
nvmlDeviceGetCurrPcieLinkGeneration.....	101
nvmlDeviceGetCurrPcieLinkWidth.....	102
nvmlDeviceGetPcieThroughput.....	102
nvmlDeviceGetPcieReplayCounter.....	103
nvmlDeviceGetClockInfo.....	104
nvmlDeviceGetMaxClockInfo.....	104
nvmlDeviceGetGpcClkVfOffset.....	105
nvmlDeviceGetApplicationsClock.....	106
nvmlDeviceGetDefaultApplicationsClock.....	106
nvmlDeviceGetClock.....	107
nvmlDeviceGetMaxCustomerBoostClock.....	108
nvmlDeviceGetSupportedMemoryClocks.....	109
nvmlDeviceGetSupportedGraphicsClocks.....	110

nvmlDeviceGetAutoBoostedClocksEnabled.....	111
nvmlDeviceGetFanSpeed.....	112
nvmlDeviceGetFanSpeed_v2.....	112
nvmlDeviceGetTargetFanSpeed.....	113
nvmlDeviceGetMinMaxFanSpeed.....	114
nvmlDeviceGetFanControlPolicy_v2.....	115
nvmlDeviceGetNumFans.....	115
nvmlDeviceGetTemperature.....	116
nvmlDeviceGetTemperatureThreshold.....	116
nvmlDeviceGetThermalSettings.....	118
nvmlDeviceGetPerformanceState.....	118
nvmlDeviceGetCurrentClocksEventReasons.....	119
nvmlDeviceGetCurrentClocksThrottleReasons.....	120
nvmlDeviceGetSupportedClocksEventReasons.....	120
nvmlDeviceGetSupportedClocksThrottleReasons.....	121
nvmlDeviceGetPowerState.....	121
nvmlDeviceGetDynamicPstatesInfo.....	122
nvmlDeviceGetMemClkVfOffset.....	122
nvmlDeviceGetMinMaxClockOfPState.....	123
nvmlDeviceGetSupportedPerformanceStates.....	123
nvmlDeviceGetGpcClkMinMaxVfOffset.....	124
nvmlDeviceGetMemClkMinMaxVfOffset.....	125
nvmlDeviceGetPowerManagementMode.....	126
nvmlDeviceGetPowerManagementLimit.....	127
nvmlDeviceGetPowerManagementLimitConstraints.....	127
nvmlDeviceGetPowerManagementDefaultLimit.....	128
nvmlDeviceGetPowerUsage.....	129
nvmlDeviceGetTotalEnergyConsumption.....	130
nvmlDeviceGetEnforcedPowerLimit.....	130
nvmlDeviceGetGpuOperationMode.....	131
nvmlDeviceGetMemoryInfo.....	132
nvmlDeviceGetComputeMode.....	133
nvmlDeviceGetCudaComputeCapability.....	134
nvmlDeviceGetEccMode.....	134
nvmlDeviceGetDefaultEccMode.....	135
nvmlDeviceGetBoardId.....	136
nvmlDeviceGetMultiGpuBoard.....	137
nvmlDeviceGetTotalEccErrors.....	137
nvmlDeviceGetDetailedEccErrors.....	138
nvmlDeviceGetMemoryErrorCounter.....	140
nvmlDeviceGetUtilizationRates.....	141
nvmlDeviceGetEncoderUtilization.....	142
nvmlDeviceGetEncoderCapacity.....	142

nvmlDeviceGetEncoderStats.....	143
nvmlDeviceGetEncoderSessions.....	144
nvmlDeviceGetDecoderUtilization.....	145
nvmlDeviceGetJpgUtilization.....	146
nvmlDeviceGetOfaUtilization.....	146
nvmlDeviceGetFBCStats.....	147
nvmlDeviceGetFBCSessions.....	148
nvmlDeviceGetDriverModel.....	149
nvmlDeviceGetVbiosVersion.....	150
nvmlDeviceGetBridgeChipInfo.....	150
nvmlDeviceGetComputeRunningProcesses_v3.....	151
nvmlDeviceGetGraphicsRunningProcesses_v3.....	152
nvmlDeviceGetMPSComputeRunningProcesses_v3.....	154
nvmlDeviceGetRunningProcessDetailList.....	155
nvmlDeviceOnSameBoard.....	156
nvmlDeviceGetAPIRestriction.....	157
nvmlDeviceGetSamples.....	158
nvmlDeviceGetBAR1MemoryInfo.....	159
nvmlDeviceGetViolationStatus.....	160
nvmlDeviceGetIrqNum.....	161
nvmlDeviceGetNumGpuCores.....	161
nvmlDeviceGetPowerSource.....	162
nvmlDeviceGetMemoryBusWidth.....	163
nvmlDeviceGetPcieLinkMaxSpeed.....	163
nvmlDeviceGetPcieSpeed.....	164
nvmlDeviceGetAdaptiveClockInfoStatus.....	164
nvmlDeviceGetBusType.....	165
nvmlDeviceGetGpuFabricInfo.....	165
nvmlDeviceGetGpuFabricInfoV.....	166
nvmlSystemGetConfComputeCapabilities.....	167
nvmlSystemGetConfComputeState.....	167
nvmlDeviceGetConfComputeMemSizeInfo.....	168
nvmlSystemGetConfComputeGpusReadyState.....	168
nvmlDeviceGetConfComputeProtectedMemoryUsage.....	169
nvmlDeviceGetConfComputeGpuCertificate.....	169
nvmlDeviceGetConfComputeGpuAttestationReport.....	170
nvmlSystemGetConfComputeKeyRotationThresholdInfo.....	170
nvmlSystemGetConfComputeSettings.....	171
nvmlDeviceGetGspFirmwareVersion.....	172
nvmlDeviceGetGspFirmwareMode.....	172
nvmlDeviceGetRetiredPages.....	173
nvmlDeviceGetRetiredPages_v2.....	174
nvmlDeviceGetRetiredPagesPendingStatus.....	175

nvmlDeviceGetRemappedRows.....	176
nvmlDeviceGetRowRemapperHistogram.....	177
nvmlDeviceGetArchitecture.....	177
nvmlDeviceGetClkMonStatus.....	178
nvmlDeviceGetProcessUtilization.....	178
nvmlDeviceGetProcessesUtilizationInfo.....	180
4.15.1. CPU and Memory Affinity.....	181
nvmlDeviceGetMemoryAffinity.....	181
nvmlDeviceGetCpuAffinityWithinScope.....	182
nvmlDeviceGetCpuAffinity.....	183
nvmlDeviceSetCpuAffinity.....	184
nvmlDeviceClearCpuAffinity.....	185
nvmlDeviceGetNumaNodId.....	185
NVML_AFFINITY_SCOPE_NODE.....	186
NVML_AFFINITY_SCOPE_SOCKET.....	186
4.16. Unit Commands.....	186
nvmlUnitSetLedState.....	186
4.17. Device Commands.....	187
nvmlDeviceSetPersistenceMode.....	187
nvmlDeviceSetComputeMode.....	188
nvmlDeviceSetEccMode.....	189
nvmlDeviceClearEccErrorCounts.....	190
nvmlDeviceSetDriverModel.....	191
nvmlDeviceSetGpuLockedClocks.....	192
nvmlDeviceResetGpuLockedClocks.....	193
nvmlDeviceSetMemoryLockedClocks.....	194
nvmlDeviceResetMemoryLockedClocks.....	195
nvmlDeviceSetApplicationsClocks.....	196
nvmlDeviceResetApplicationsClocks.....	197
nvmlDeviceSetAutoBoostedClocksEnabled.....	198
nvmlDeviceSetDefaultAutoBoostedClocksEnabled.....	199
nvmlDeviceSetDefaultFanSpeed_v2.....	200
nvmlDeviceSetFanControlPolicy.....	200
nvmlDeviceSetTemperatureThreshold.....	201
nvmlDeviceSetPowerManagementLimit.....	201
nvmlDeviceSetGpuOperationMode.....	202
nvmlDeviceSetAPIRestriction.....	203
nvmlDeviceSetFanSpeed_v2.....	204
nvmlDeviceSetGpcClkVfOffset.....	205
nvmlDeviceSetMemClkVfOffset.....	205
nvmlDeviceSetConfComputeUnprotectedMemSize.....	206
nvmlSystemSetConfComputeGpusReadyState.....	206
nvmlSystemSetConfComputeKeyRotationThresholdInfo.....	207

4.18. NvLink Methods.....	208
nvmldDeviceGetNvLinkState.....	208
nvmldDeviceGetNvLinkVersion.....	208
nvmldDeviceGetNvLinkCapability.....	209
nvmldDeviceGetNvLinkRemotePciInfo_v2.....	210
nvmldDeviceGetNvLinkErrorCounter.....	210
nvmldDeviceResetNvLinkErrorCounters.....	211
nvmldDeviceSetNvLinkUtilizationControl.....	212
nvmldDeviceGetNvLinkUtilizationControl.....	213
nvmldDeviceGetNvLinkUtilizationCounter.....	213
nvmldDeviceFreezeNvLinkUtilizationCounter.....	214
nvmldDeviceResetNvLinkUtilizationCounter.....	215
nvmldDeviceGetNvLinkRemoteDeviceType.....	216
4.19. Event Handling Methods.....	217
nvmldEventData_t.....	217
Event Types.....	217
nvmldEventSet_t.....	217
nvmldEventSetCreate.....	217
nvmldDeviceRegisterEvents.....	218
nvmldDeviceGetSupportedEventTypes.....	219
nvmldEventSetWait_v2.....	220
nvmldEventSetFree.....	221
4.19.1. Event Types.....	221
nvmldEventTypeSingleBitEccError.....	222
nvmldEventTypeDoubleBitEccError.....	222
nvmldEventTypePState.....	222
nvmldEventTypeXidCriticalError.....	222
nvmldEventTypeClock.....	222
nvmldEventTypePowerSourceChange.....	222
nvmldEventMigConfigChange.....	222
nvmldEventTypeNone.....	222
nvmldEventTypeAll.....	223
4.20. Drain states.....	223
nvmldDeviceModifyDrainState.....	223
nvmldDeviceQueryDrainState.....	224
nvmldDeviceRemoveGpu_v2.....	224
nvmldDeviceDiscoverGpus.....	225
4.21. Field Value Queries.....	226
nvmldDeviceGetFieldValues.....	226
nvmldDeviceClearFieldValues.....	227
4.22. Enums, Constants and Structs.....	227
4.23. vGPU APIs.....	227
nvmldDeviceGetVirtualizationMode.....	228

nvmlDeviceGetHostVgpuMode.....	228
nvmlDeviceSetVirtualizationMode.....	229
nvmlDeviceGetVgpuHeterogeneousMode.....	230
nvmlDeviceSetVgpuHeterogeneousMode.....	231
nvmlVgpuInstanceGetPlacementId.....	232
nvmlDeviceGetVgpuTypeSupportedPlacements.....	232
nvmlDeviceGetVgpuTypeCreatablePlacements.....	233
nvmlVgpuTypeGetGspHeapSize.....	234
nvmlVgpuTypeGetFbReservation.....	235
nvmlDeviceSetVgpuCapabilities.....	235
nvmlDeviceGetGridLicensableFeatures_v4.....	236
4.24. vGPU Management.....	236
nvmlGetVgpuDriverCapabilities.....	237
nvmlDeviceGetVgpuCapabilities.....	237
nvmlDeviceGetSupportedVgpus.....	238
nvmlDeviceGetCreatableVgpus.....	239
nvmlVgpuTypeGetClass.....	240
nvmlVgpuTypeGetName.....	241
nvmlVgpuTypeGetGpuInstanceProfileId.....	241
nvmlVgpuTypeGetDeviceId.....	242
nvmlVgpuTypeGetFramebufferSize.....	243
nvmlVgpuTypeGetNumDisplayHeads.....	243
nvmlVgpuTypeGetResolution.....	244
nvmlVgpuTypeGetLicense.....	245
nvmlVgpuTypeGetFrameRateLimit.....	246
nvmlVgpuTypeGetMaxInstances.....	246
nvmlVgpuTypeGetMaxInstancesPerVm.....	247
nvmlDeviceGetActiveVgpus.....	248
nvmlVgpuInstanceGetVmId.....	249
nvmlVgpuInstanceGetUUID.....	250
nvmlVgpuInstanceGetVmDriverVersion.....	250
nvmlVgpuInstanceGetFbUsage.....	251
nvmlVgpuInstanceGetLicenseStatus.....	252
nvmlVgpuInstanceGetType.....	253
nvmlVgpuInstanceGetFrameRateLimit.....	253
nvmlVgpuInstanceGetEccMode.....	254
nvmlVgpuInstanceGetEncoderCapacity.....	255
nvmlVgpuInstanceSetEncoderCapacity.....	255
nvmlVgpuInstanceGetEncoderStats.....	256
nvmlVgpuInstanceGetEncoderSessions.....	257
nvmlVgpuInstanceGetFBCStats.....	258
nvmlVgpuInstanceGetFBCSessions.....	258
nvmlVgpuInstanceGetGpuInstanceId.....	259

nvmiVgpuInstanceGetGpuPcild.....	260
nvmiVgpuTypeGetCapabilities.....	261
nvmiVgpuInstanceGetMdevUUID.....	261
4.25. vGPU Migration.....	262
nvmiVgpuVersion_t.....	263
nvmiVgpuMetadata_t.....	263
nvmiVgpuPgpuMetadata_t.....	263
nvmiVgpuPgpuCompatibility_t.....	263
nvmiVgpuVmCompatibility_t.....	263
nvmiVgpuPgpuCompatibilityLimitCode_t.....	263
nvmiVgpuInstanceGetMetadata.....	264
nvmiDeviceGetVgpuMetadata.....	265
nvmiGetVgpuCompatibility.....	265
nvmiDeviceGetPgpuMetadataString.....	266
nvmiDeviceGetVgpuSchedulerLog.....	267
nvmiDeviceGetVgpuSchedulerState.....	268
nvmiDeviceGetVgpuSchedulerCapabilities.....	269
nvmiDeviceSetVgpuSchedulerState.....	269
nvmiGetVgpuVersion.....	270
nvmiSetVgpuVersion.....	271
4.26. vGPU Utilization and Accounting.....	272
nvmiDeviceGetVgpuUtilization.....	272
nvmiDeviceGetVgpuInstancesUtilizationInfo.....	274
nvmiDeviceGetVgpuProcessUtilization.....	275
nvmiDeviceGetVgpuProcessesUtilizationInfo.....	277
nvmiVgpuInstanceGetAccountingMode.....	278
nvmiVgpuInstanceGetAccountingPids.....	279
nvmiVgpuInstanceGetAccountingStats.....	280
nvmiVgpuInstanceClearAccountingPids.....	281
nvmiVgpuInstanceGetLicenseInfo_v2.....	282
4.27. Excluded GPU Queries.....	282
nvmiExcludedDeviceInfo_t.....	282
nvmiGetExcludedDeviceCount.....	282
nvmiGetExcludedDeviceInfoByIndex.....	283
4.28. Multi Instance GPU Management.....	283
nvmiGpuInstancePlacement_t.....	284
nvmiGpuInstanceProfileInfo_t.....	284
nvmiGpuInstanceProfileInfo_v2_t.....	284
nvmiGpuInstanceProfileInfo_v3_t.....	284
nvmiComputeInstanceProfileInfo_t.....	284
nvmiComputeInstanceProfileInfo_v2_t.....	284
nvmiComputeInstanceProfileInfo_v3_t.....	284
nvmiDeviceSetMigMode.....	284



nvmlDeviceGetMigMode.....	285
nvmlDeviceGetGpuInstanceProfileInfo.....	286
nvmlDeviceGetGpuInstanceProfileInfoV.....	287
nvmlDeviceGetGpuInstancePossiblePlacements_v2.....	287
nvmlDeviceGetGpuInstanceRemainingCapacity.....	288
nvmlDeviceCreateGpuInstance.....	289
nvmlDeviceCreateGpuInstanceWithPlacement.....	290
nvmlGpuInstanceDestroy.....	291
nvmlDeviceGetGpuInstances.....	292
nvmlDeviceGetGpuInstanceById.....	293
nvmlGpuInstanceGetInfo.....	293
nvmlGpuInstanceGetComputeInstanceProfileInfo.....	294
nvmlGpuInstanceGetComputeInstanceProfileInfoV.....	295
nvmlGpuInstanceGetComputeInstanceRemainingCapacity.....	296
nvmlGpuInstanceGetComputeInstancePossiblePlacements.....	296
nvmlGpuInstanceCreateComputeInstance.....	298
nvmlGpuInstanceCreateComputeInstanceWithPlacement.....	298
nvmlComputeInstanceDestroy.....	300
nvmlGpuInstanceGetComputeInstances.....	300
nvmlGpuInstanceGetComputeInstanceById.....	301
nvmlComputeInstanceGetInfo_v2.....	302
nvmlDeviceMigDeviceHandle.....	302
nvmlDeviceGetGpuInstanceId.....	303
nvmlDeviceGetComputeInstanceId.....	304
nvmlDeviceGetMaxMigDeviceCount.....	304
nvmlDeviceGetMigDeviceHandleByIndex.....	305
nvmlDeviceGetDeviceHandleFromMigDeviceHandle.....	306
NVML_DEVICE_MIG_DISABLE.....	306
NVML_DEVICE_MIG_ENABLE.....	306
NVML_GPU_INSTANCE_PROFILE_1_SLICE.....	306
NVML_GPU_INSTANCE_PROFILE_CAPS_P2P.....	306
nvmlGpuInstanceProfileInfo_v2.....	307
nvmlGpuInstanceProfileInfo_v3.....	307
NVML_COMPUTE_INSTANCE_PROFILE_1_SLICE.....	307
NVML_COMPUTE_INSTANCE_ENGINE_PROFILE_SHARED.....	307
nvmlComputeInstanceProfileInfo_v2.....	307
nvmlComputeInstanceProfileInfo_v3.....	307
4.29. NVML GPM.....	307
GPM Enums.....	308
GPM Structs.....	308
GPM Functions.....	308
4.29.1. GPM Enums.....	308
nvmlGpmMetricId_t.....	308

4.29.2. GPM Structs.....	311
nvmGpmMetric_t.....	312
nvmGpmMetricsGet_t.....	312
nvmGpmSupport_t.....	312
nvmGpmSample_t.....	312
4.29.3. GPM Functions.....	312
nvmGpmMetricsGet.....	312
nvmGpmSampleFree.....	312
nvmGpmSampleAlloc.....	313
nvmGpmSampleGet.....	313
nvmGpmMigSampleGet.....	314
nvmGpmQueryDeviceSupport.....	314
nvmGpmQueryIfStreamingEnabled.....	315
nvmGpmSetStreamingEnabled.....	315
4.30. VirtualGPU.....	316
vGPU Enums.....	316
vGPU Constants.....	316
vGPU Structs.....	316
4.30.1. vGPU Enums.....	316
nvmGpuVirtualizationMode_t.....	316
nvmHostVgpuMode_t.....	316
nvmVgpuVmlDType_t.....	317
nvmVgpuGuestInfoState_t.....	317
nvmGridLicenseFeatureCode_t.....	317
nvmVgpuCapability_t.....	318
nvmVgpuDriverCapability_t.....	318
nvmDeviceVgpuCapability_t.....	318
NVML_GRID_LICENSE_EXPIRY_NOT_AVAILABLE.....	319
NVML_GRID_LICENSE_EXPIRY_INVALID.....	319
NVML_GRID_LICENSE_EXPIRY_VALID.....	319
NVML_GRID_LICENSE_EXPIRY_NOT_APPLICABLE.....	319
NVML_GRID_LICENSE_EXPIRY_PERMANENT.....	319
4.30.2. vGPU Constants.....	319
NVML_GRID_LICENSE_BUFFER_SIZE.....	319
NVML_VGPU_VIRTUALIZATION_CAP_MIGRATION.....	319
NVML_VGPU_PGPU_VIRTUALIZATION_CAP_MIGRATION.....	319
4.30.3. vGPU Structs.....	320
nvmVgpuHeterogeneousMode_v1_t.....	321
nvmVgpuPlacementId_v1_t.....	321
nvmVgpuPlacementList_v1_t.....	321
nvmVgpuInstanceUtilizationSample_t.....	321
nvmVgpuInstanceUtilizationInfo_v1_t.....	321
nvmVgpuInstancesUtilizationInfo_v1_t.....	321

nvmlVgpuProcessUtilizationSample_t.....	321
nvmlVgpuProcessUtilizationInfo_v1_t.....	321
nvmlVgpuProcessesUtilizationInfo_v1_t.....	321
nvmlVgpuSchedulerParams_t.....	321
nvmlVgpuSchedulerLogEntry_t.....	321
nvmlVgpuSchedulerLog_t.....	321
nvmlVgpuSchedulerGetState_t.....	321
nvmlVgpuSchedulerSetParams_t.....	321
nvmlVgpuSchedulerSetState_t.....	321
nvmlVgpuSchedulerCapabilities_t.....	321
nvmlVgpuLicenseExpiry_t.....	321
nvmlProcessUtilizationSample_t.....	321
nvmlProcessUtilizationInfo_v1_t.....	321
nvmlProcessesUtilizationInfo_v1_t.....	321
nvmlGridLicenseExpiry_t.....	321
nvmlGridLicensableFeature_t.....	322
nvmlGridLicensableFeatures_t.....	322
nvmlEccSramErrorStatus_v1_t.....	322
NVML_VGPU_SCHEDULER_POLICY_UNKNOWN.....	322
NVML_GRID_LICENSE_STATE_UNKNOWN.....	322
NVML_GRID_LICENSE_STATE_UNINITIALIZED.....	322
NVML_GRID_LICENSE_STATE_UNLICENSED_UNRESTRICTED.....	322
NVML_GRID_LICENSE_STATE_UNLICENSED_RESTRICTED.....	322
NVML_GRID_LICENSE_STATE_UNLICENSED.....	322
NVML_GRID_LICENSE_STATE_LICENSED.....	322
NVML_GSP_FIRMWARE_VERSION_BUF_SIZE.....	322
NVML_DEVICE_ARCH_KEPLER.....	322
NVML_BUS_TYPE_UNKNOWN.....	322
NVML_FAN_POLICY_TEMPERATURE_CONTINUOUS_SW.....	323
NVML_POWER_SOURCE_AC.....	323
4.31. NvmlClocksEventReasons.....	323
nvmlClocksEventReasonGpudle.....	323
nvmlClocksEventReasonApplicationsClocksSetting.....	323
nvmlClocksThrottleReasonUserDefinedClocks.....	323
nvmlClocksEventReasonSwPowerCap.....	323
nvmlClocksThrottleReasonHwSlowdown.....	324
nvmlClocksEventReasonSyncBoost.....	324
nvmlClocksEventReasonSwThermalSlowdown.....	324
nvmlClocksThrottleReasonHwThermalSlowdown.....	325
nvmlClocksThrottleReasonHwPowerBrakeSlowdown.....	325
nvmlClocksEventReasonDisplayClockSetting.....	325
nvmlClocksEventReasonNone.....	326
nvmlClocksEventReasonAll.....	326

nvmcClocksThrottleReasonGpuldle.....	326
nvmcClocksThrottleReasonApplicationsClocksSetting.....	326
nvmcClocksThrottleReasonSyncBoost.....	326
nvmcClocksThrottleReasonSwPowerCap.....	327
nvmcClocksThrottleReasonSwThermalSlowdown.....	327
nvmcClocksThrottleReasonDisplayClockSetting.....	327
nvmcClocksThrottleReasonNone.....	327
nvmcClocksThrottleReasonAll.....	327
<b>Chapter 5. Data Structures.....</b>	<b>328</b>
nvmAccountingStats_t.....	330
gpuUtilization.....	330
memoryUtilization.....	330
maxMemoryUsage.....	330
time.....	330
startTime.....	331
isRunning.....	331
reserved.....	331
nvmBAR1Memory_t.....	331
bar1Total.....	331
bar1Free.....	331
bar1Used.....	331
nvmBridgeChipHierarchy_t.....	331
bridgeCount.....	332
bridgeChipInfo.....	332
nvmBridgeChipInfo_t.....	332
type.....	332
fwVersion.....	332
nvmC2cModelInfo_v1_t.....	332
nvmClkMonFaultInfo_t.....	332
clkApiDomain.....	332
clkDomainFaultMask.....	332
nvmClkMonStatus_t.....	333
bGlobalStatus.....	333
clkMonListSize.....	333
clkMonList.....	333
nvmComputeInstanceProfileInfo_t.....	333
id.....	334
sliceCount.....	334
instanceCount.....	334
multiprocessorCount.....	334
sharedCopyEngineCount.....	334
sharedDecoderCount.....	334
sharedEncoderCount.....	334

sharedJpegCount.....	334
sharedOfaCount.....	334
nvmlComputeInstanceProfileInfo_v2_t.....	335
version.....	336
id.....	336
sliceCount.....	336
instanceCount.....	336
multiprocessorCount.....	336
sharedCopyEngineCount.....	336
sharedDecoderCount.....	336
sharedEncoderCount.....	336
sharedJpegCount.....	336
sharedOfaCount.....	337
name.....	337
nvmlComputeInstanceProfileInfo_v3_t.....	337
version.....	338
id.....	338
sliceCount.....	338
instanceCount.....	338
multiprocessorCount.....	338
sharedCopyEngineCount.....	338
sharedDecoderCount.....	338
sharedEncoderCount.....	338
sharedJpegCount.....	338
sharedOfaCount.....	339
name.....	339
capabilities.....	339
nvmlConfComputeMemSizeInfo_t.....	339
nvmlEccErrorCounts_t.....	339
l1Cache.....	340
l2Cache.....	340
deviceMemory.....	340
registerFile.....	340
nvmlEccSramErrorStatus_v1_t.....	340
version.....	341
aggregateUncParity.....	341
aggregateUncSecDed.....	341
aggregateCor.....	341
volatileUncParity.....	341
volatileUncSecDed.....	341
volatileCor.....	341
aggregateUncBucketL2.....	341
aggregateUncBucketSm.....	341

aggregateUncBucketPcie.....	342
aggregateUncBucketMcu.....	342
aggregateUncBucketOther.....	342
bThresholdExceeded.....	342
nvmEncoderSessionInfo_t.....	342
sessionId.....	343
pid.....	343
vgpuInstance.....	343
codecType.....	343
hResolution.....	343
vResolution.....	343
averageFps.....	343
averageLatency.....	343
nvmEventData_t.....	343
device.....	344
eventType.....	344
eventData.....	344
gpuInstanceId.....	344
computeInstanceId.....	344
nvmExcludedDeviceInfo_t.....	344
pciInfo.....	344
uuid.....	344
nvmFBCSessionInfo_t.....	344
sessionId.....	345
pid.....	345
vgpuInstance.....	345
displayOrdinal.....	345
sessionType.....	345
sessionFlags.....	345
hMaxResolution.....	345
vMaxResolution.....	345
hResolution.....	345
vResolution.....	345
averageFPS.....	345
averageLatency.....	345
nvmFBCStats_t.....	346
sessionsCount.....	346
averageFPS.....	346
averageLatency.....	346
nvmFieldValue_t.....	346
fieldId.....	347
scopeId.....	347
timestamp.....	347

latencyUsec.....	347
valueType.....	347
nvmlReturn.....	347
value.....	347
nvmlGpmMetric_t.....	347
metricId.....	348
nvmlReturn.....	348
value.....	348
metricInfo.....	348
nvmlGpmMetricsGet_t.....	348
version.....	348
numMetrics.....	348
sample1.....	348
sample2.....	348
metrics.....	348
nvmlGpmSupport_t.....	348
version.....	349
isSupportedDevice.....	349
nvmlGpuFabricInfo_v2_t.....	349
version.....	349
clusterUuid.....	349
status.....	349
cliqueId.....	349
state.....	349
healthMask.....	349
nvmlGpuInstancePlacement_t.....	349
start.....	350
size.....	350
nvmlGpuInstanceProfileInfo_t.....	350
id.....	351
isP2pSupported.....	351
sliceCount.....	351
instanceCount.....	351
multiprocessorCount.....	351
copyEngineCount.....	351
decoderCount.....	351
encoderCount.....	351
jpegCount.....	351
ofaCount.....	351
memorySizeMB.....	352
nvmlGpuInstanceProfileInfo_v2_t.....	352
version.....	353
id.....	353

isP2pSupported.....	353
sliceCount.....	353
instanceCount.....	353
multiprocessorCount.....	353
copyEngineCount.....	353
decoderCount.....	353
encoderCount.....	353
jpegCount.....	353
ofaCount.....	354
memorySizeMB.....	354
name.....	354
nvmiGpuInstanceProfileInfo_v3_t.....	354
version.....	355
id.....	355
sliceCount.....	355
instanceCount.....	355
multiprocessorCount.....	355
copyEngineCount.....	355
decoderCount.....	355
encoderCount.....	355
jpegCount.....	355
ofaCount.....	355
memorySizeMB.....	356
name.....	356
capabilities.....	356
nvmiGridLicensableFeature_t.....	356
featureCode.....	357
featureState.....	357
licenseInfo.....	357
productName.....	357
featureEnabled.....	357
licenseExpiry.....	357
nvmiGridLicensableFeatures_t.....	357
isGridLicenseSupported.....	358
licensableFeaturesCount.....	358
gridLicensableFeatures.....	358
nvmiGridLicenseExpiry_t.....	358
year.....	359
month.....	359
day.....	359
hour.....	359
min.....	359
sec.....	359



status.....	359
nvmlHwbcEntry_t.....	359
nvmlLedState_t.....	359
cause.....	360
color.....	360
nvmlMemory_t.....	360
total.....	360
free.....	360
used.....	360
nvmlMemory_v2_t.....	360
version.....	361
total.....	361
reserved.....	361
free.....	361
used.....	361
nvmlNvLinkUtilizationControl_t.....	361
nvmlPciInfo_t.....	361
busIdLegacy.....	362
domain.....	362
bus.....	362
device.....	362
pciDeviceId.....	362
pciSubSystemId.....	362
busId.....	362
nvmlPciInfoExt_v1_t.....	362
version.....	363
domain.....	363
bus.....	363
device.....	363
pciDeviceId.....	363
pciSubSystemId.....	363
baseClass.....	363
subClass.....	363
busId.....	363
nvmlProcessDetail_v1_t.....	363
pid.....	364
usedGpuMemory.....	364
gpuInstanceId.....	364
computeInstanceId.....	364
usedGpuCcProtectedMemory.....	364
nvmlProcessDetailList_v1_t.....	364
version.....	365
mode.....	365

numProcArrayEntries.....	365
procArray.....	365
nvmlProcessesUtilizationInfo_v1_t.....	365
version.....	366
processSamplesCount.....	366
lastSeenTimeStamp.....	366
procUtilArray.....	366
nvmlProcessInfo_t.....	366
pid.....	366
usedGpuMemory.....	366
gpuInstanceId.....	367
computeInstanceId.....	367
nvmlProcessInfo_v1_t.....	367
pid.....	367
usedGpuMemory.....	367
nvmlProcessUtilizationInfo_v1_t.....	367
timeStamp.....	368
pid.....	368
smUtil.....	368
memUtil.....	368
encUtil.....	368
decUtil.....	368
jpgUtil.....	368
ofaUtil.....	368
nvmlProcessUtilizationSample_t.....	368
pid.....	369
timeStamp.....	369
smUtil.....	369
memUtil.....	369
encUtil.....	369
decUtil.....	369
nvmlPSUInfo_t.....	369
state.....	370
current.....	370
voltage.....	370
power.....	370
nvmlRowRemapperHistogramValues_t.....	370
nvmlSample_t.....	370
timeStamp.....	370
sampleValue.....	370
nvmlSystemConfComputeSettings_v1_t.....	370
nvmlUnitFanInfo_t.....	371
speed.....	371

state.....	371
nvmUnitFanSpeeds_t.....	371
fans.....	371
count.....	371
nvmUnitInfo_t.....	371
name.....	372
id.....	372
serial.....	372
firmwareVersion.....	372
nvmUtilization_t.....	372
gpu.....	372
memory.....	372
nvmValue_t.....	372
dVal.....	373
siVal.....	373
uiVal.....	373
ulVal.....	373
ullVal.....	373
sllVal.....	373
nvmVgpuHeterogeneousMode_v1_t.....	373
version.....	373
mode.....	373
nvmVgpuInstancesUtilizationInfo_v1_t.....	373
version.....	374
sampleValType.....	374
vgpuInstanceCount.....	374
lastSeenTimeStamp.....	374
vgpuUtilArray.....	374
nvmVgpuInstanceUtilizationInfo_v1_t.....	374
timeStamp.....	375
vgpuInstance.....	375
smUtil.....	375
memUtil.....	375
encUtil.....	375
decUtil.....	375
jpgUtil.....	375
ofaUtil.....	375
nvmVgpuInstanceUtilizationSample_t.....	375
vgpuInstance.....	376
timeStamp.....	376
smUtil.....	376
memUtil.....	376
encUtil.....	376

decUtil.....	376
nvmLVgpuLicenseExpiry_t.....	376
year.....	377
month.....	377
day.....	377
hour.....	377
min.....	377
sec.....	377
status.....	377
nvmLVgpuMetadata_t.....	377
version.....	378
revision.....	378
guestInfoState.....	378
guestDriverVersion.....	378
hostDriverVersion.....	378
reserved.....	378
vgpuVirtualizationCaps.....	378
guestVgpuVersion.....	378
opaqueDataSize.....	378
opaqueData.....	378
nvmLVgpuPgpuCompatibility_t.....	378
vgpuVmCompatibility.....	379
compatibilityLimitCode.....	379
nvmLVgpuPgpuMetadata_t.....	379
version.....	380
revision.....	380
hostDriverVersion.....	380
pgpuVirtualizationCaps.....	380
reserved.....	380
hostSupportedVgpuRange.....	380
opaqueDataSize.....	380
opaqueData.....	380
nvmLVgpuPlacementId_v1_t.....	380
version.....	381
placementId.....	381
nvmLVgpuPlacementList_v1_t.....	381
version.....	381
placementSize.....	381
count.....	381
placementIds.....	381
nvmLVgpuProcessesUtilizationInfo_v1_t.....	381
version.....	382
vgpuProcessCount.....	382

lastSeenTimeStamp.....	382
vgpuProcUtilArray.....	382
nvmlVgpuProcessUtilizationInfo_v1_t.....	382
processName.....	383
timeStamp.....	383
vgpuInstance.....	383
pid.....	383
smUtil.....	383
memUtil.....	383
encUtil.....	383
decUtil.....	383
jpgUtil.....	383
ofaUtil.....	383
nvmlVgpuProcessUtilizationSample_t.....	384
vgpuInstance.....	385
pid.....	385
processName.....	385
timeStamp.....	385
smUtil.....	385
memUtil.....	385
encUtil.....	385
decUtil.....	385
nvmlVgpuSchedulerCapabilities_t.....	385
supportedSchedulers.....	386
maxTimeslice.....	386
minTimeslice.....	386
isArrModeSupported.....	386
maxFrequencyForARR.....	386
minFrequencyForARR.....	386
maxAvgFactorForARR.....	386
minAvgFactorForARR.....	386
nvmlVgpuSchedulerGetState_t.....	386
schedulerPolicy.....	387
arrMode.....	387
nvmlVgpuSchedulerLog_t.....	387
engineId.....	387
schedulerPolicy.....	387
arrMode.....	387
entriesCount.....	387
nvmlVgpuSchedulerLogEntry_t.....	387
timestamp.....	388
timeRunTotal.....	388
timeRun.....	388

swRunlistId.....	388
targetTimeSlice.....	388
cumulativePreemptionTime.....	388
nvmiVgpuSchedulerParams_t.....	388
avgFactor.....	389
timeslice.....	389
nvmiVgpuSchedulerSetParams_t.....	389
avgFactor.....	389
frequency.....	389
timeslice.....	389
nvmiVgpuSchedulerSetState_t.....	389
schedulerPolicy.....	390
enableARRMode.....	390
nvmiVgpuVersion_t.....	390
minVersion.....	390
maxVersion.....	390
nvmiViolationTime_t.....	390
referenceTime.....	390
violationTime.....	390
<b>Chapter 6. Data Fields.....</b>	<b>391</b>
<b>Chapter 7. Deprecated List.....</b>	<b>406</b>

# Chapter 1.

## NVML API REFERENCE

The NVIDIA Management Library (*NVML*) is a C-based programmatic interface for monitoring and managing various states within NVIDIA Tesla™ GPUs. It is intended to be a platform for building 3rd party applications, and is also the underlying library for the NVIDIA-supported `nvidia-smi` tool. NVML is thread-safe so it is safe to make simultaneous NVML calls from multiple threads.

### API Documentation

Supported OS platforms:

- ▶ Windows 64-bit: Windows Server 2019, Windows Server 2016, Windows Server 2012 R2, Windows 10
- ▶ Linux: 64-bit

Supported products:

- ▶ Full Support
  - ▶ NVIDIA Tesla Line:
    - ▶ A100, A40, A30, A16, A10
    - ▶ H100
    - ▶ T4
    - ▶ V100
    - ▶ P100, P40, P4, P6
    - ▶ M60, M40, M6, M4
    - ▶ K80, K520
  - ▶ NVIDIA Quadro Line:
    - ▶ RTX 8000, RTX 6000, RTX 5000, RTX 4000, RTX 3000
    - ▶ GV100, GP100, P6000, P5200, P5000, P4000, P2200, P2000, P1000, P620, P600, P400
    - ▶ M6000 24GB, M6000, M5000, M4000, M2000

- ▶ K6000, K5200, K5000, K4000, K4200, K2200, K2000, K2000D, K1200, K620, K600, K420, 410
- ▶ NVIDIA GeForce Line:
  - ▶ none
- ▶ Limited Support
  - ▶ NVIDIA Tesla Line: All other current and previous generation Tesla-branded parts
  - ▶ NVIDIA Quadro Line: All other current and previous generation Quadro-branded parts
  - ▶ NVIDIA GeForce Line: All current and previous generation GeForce-branded parts

The NVML library can be found at the following locations on Windows

- ▶ Standard driver install: %ProgramW6432%\**"NVIDIA Corporation"**\NVSMI\
- ▶ DCH driver install: \Windows\System32

Note that these libraries will not be added to the path on Windows. To dynamically link to NVML, add this path to the PATH environmental variable. To dynamically load NVML, call **LoadLibrary** with this path.

On Linux the NVML library is named "**libnvidia-ml.so**" and can be found on the standard library path. To link against the NVML library add the **-lnvidia-ml** flag to your linker command.

The NVML API is divided into five categories:

- ▶ Support Methods:
  - ▶ Initialization and Cleanup
- ▶ Query Methods:
  - ▶ System Queries
  - ▶ Device Queries
  - ▶ Unit Queries
- ▶ Control Methods:
  - ▶ Unit Commands
  - ▶ Device Commands
- ▶ Event Handling Methods:
  - ▶ Event Handling Methods
- ▶ Error reporting Methods
  - ▶ Error Reporting



## Chapter 2. KNOWN ISSUES

This is a list of known NVML issues in the current driver:

- ▶ On systems where GPUs are NUMA nodes, the accuracy of FB memory utilization provided by NVML depends on the memory accounting of the operating system. This is because FB memory is managed by the operating system instead of the NVIDIA GPU driver. Typically, pages allocated from FB memory are not released even after the process terminates to enhance performance. In scenarios where the operating system is under memory pressure, it may resort to utilizing FB memory. Such actions can result in discrepancies in the accuracy of memory reporting.
- ▶ On Linux, GPU Reset can't be triggered when there is pending GPU Operation Mode (GOM) change.
- ▶ On Linux, GPU Reset may not successfully change pending ECC mode. A full reboot may be required to enable the mode change.
- ▶ **`nvmlAccountingStats`** supports only one process per GPU at a time (CUDA proxy server counts as one process).
- ▶ **`nvmlAccountingStats_t.time`** reports time and utilization values starting from `cuInit` till process termination. Future driver versions might change this behavior slightly and account process only from **`cuCtxCreate`** till **`cuCtxDestroy`**.
- ▶ On GPUs from Fermi family, current P0 clocks (reported by **`nvmlDeviceGetClockInfo`**) can differ from max clocks by a few MHz.

# Chapter 3.

## CHANGE LOG

This chapter lists changes in API and bug fixes that were introduced to the library.

### Changes between v545 and v550

The following new functionality is exposed on NVIDIA display drivers version 550 Production or later.

- ▶ Added `nvmlDeviceGetNumaNodeId` to query the NUMA node of a GPU.
- ▶ Added new GPM metric ID `NVML_GPM_METRIC_NVOFA_1_UTIL` to `nvmlGpmMetricId_t`.
- ▶ Added new field ID `NVML_FI_DEV_IS_MIG_MODE_INDEPENDENT_MIG_QUERY_CAPABLE`, to check MIG query capable device irrespective of MIG mode.
- ▶ Deprecated `NVML_P2P_CAPS_INDEX_PROP` and added `NVML_P2P_CAPS_INDEX_PCI` to reflect the same P2P capability.
- ▶ Added `nvmlDeviceGetProcessesUtilizationInfo` to retrieve the recent utilization and process ID for all running processes.
- ▶ Added new `struct nvmlProcessesUtilizationInfo_t`, which includes the new utilization of NVJPG and NVOFA.
- ▶ Added `nvmlDeviceGetVgpuInstancesUtilizationInfo` to retrieve the recent utilization for vGPU instances running on a physical GPU.
- ▶ Added `nvmlDeviceGetVgpuProcessesUtilizationInfo` to retrieve the recent utilization for processes running on vGPU instances on a physical GPU.
- ▶ Added `nvmlDeviceSetVgpuHeterogeneousMode` to enable or disable vGPU heterogeneous mode for the device.
- ▶ Added `nvmlDeviceGetVgpuHeterogeneousMode` to query the vGPU heterogeneous mode for the device.
- ▶ Added `nvmlVgpuInstanceGetPlacementId` to query placement ID of the active vGPU instance.

- ▶ Added `nvmlDeviceGetVgpuTypeSupportedPlacements` to query the supported vGPU placement IDs of a vGPU type.
- ▶ Added `nvmlDeviceGetVgpuTypeCreatablePlacements` to query the creatable vGPU placement IDs of a vGPU type.
- ▶ Added support to display confidential compute protected memory along with `fb` and `bar1` in `nvidia-smi pmon` and `dmon` commands.
- ▶ Added `nvmlDeviceGetGpuFabricInfoV` to query GPU Fabric Probe Info for the device.
- ▶ Deprecated `nvmlDeviceGetGpuFabricInfo`. This function should not be used, and will be removed in a future release. Use `nvmlDeviceGetGpuFabricInfoV` instead.
- ▶ Modified `nvmlDeviceGetGpuInstanceProfileInfo` and `nvmlDeviceGetGpuInstancePossiblePlacements` to no longer require MIG being enabled.
- ▶ Added new encoder type `NVML_ENCODER_QUERY_AV1` and `NVML_ENCODER_QUERY_UNKNOWN` to enumeration `nvmlEncoderType_t`.xxxxxxxxxxxxxxxx
- ▶ Added `nvmlSystemSetConfComputeKeyRotationThresholdInfo` to set confidential compute key rotation threshold.
- ▶ Added `nvmlSystemGetConfComputeKeyRotationThresholdInfo` to query confidential compute key rotation threshold detail.
- ▶ Added `nvmlDeviceSetVgpuCapabilities` to set the desirable vGPU capability of a device.

### Changes between v535 and v545

The following new functionality is exposed on NVIDIA display drivers version 545 Production or later.

- ▶ Added a new error code `NVML_ERROR_GPU_NOT_FOUND` to be returned if no supported GPUS are found during initialization.
- ▶ In `nvmlGpuFabricInfo_t`, `partitionId` has been renamed to `cliqId`.
- ▶ Added new versioned structs `nvmlGpuInstanceProfileInfo_v3_t` and `nvmlComputeInstanceProfileInfo_v3_t`.
- ▶ Added `nvmlDeviceGetLastBBXFlushTime` for retrieving the timestamp and duration of the latest flush of the BBX object to the inforom storage.
- ▶ Added `NVML_POWER_SCOPE_MEMORY` to report out power usage for GPU Memory.
- ▶ Added `nvmlDeviceGetPciInfoExt` which expands `nvmlDeviceGetPciInfo` to also report PCI base and sub classcodes.
- ▶ Added new struct `nvmlPciInfoExt_t`, which is used in `nvmlDeviceGetPciInfoExt`.
- ▶ Added `nvmlDeviceGetRunningProcessDetailList` API to get information about Compute, Graphics or MPS-Compute processes running on a GPU with protected memory usage info.

## Changes between v530 and v535

The following new functionality is exposed on NVIDIA display drivers version 535 Production or later.

- ▶ Added **nvmlDeviceGetSramEccErrorStatus** to query SRAM ECC error status for the device.
- ▶ Added **nvmlDeviceGetModuleId** for getting device module id.
- ▶ Updated **nvmlDeviceGetPowerSource** API to report undersized power source.
- ▶ Added **nvmlDeviceGetJpgUtilization** and **nvmlDeviceGetOfaUtilization** APIs.
- ▶ Added **nvmlSystemGetNvlinkBwMode** and **nvmlSystemSetNvlinkBwMode** APIs.
- ▶ Added **nvmlDeviceSetVgpuSchedulerState** to set the vGPU scheduler state.
- ▶ Added new field Id **NVML\_FI\_DEV\_IS\_RESETLESS\_MIG\_SUPPORTED** for device's resetless MIG capability.
- ▶ Added **nvmlDeviceGetComputeRunningProcesses\_v3** to get information about Compute processes running on a GPU.
- ▶ Added **nvmlDeviceGetGraphicsRunningProcesses\_v3** to get information about Graphics processes running on a GPU.
- ▶ Added **nvmlDeviceGetMPSComputeRunningProcesses\_v3** to get information about MPS-Compute processes running on a GPU.
- ▶ Added **nvmlDeviceGetRunningProcessDetailList** to get information about Compute, Graphics or MPS-Compute processes running on a GPU with protected memory usage info.
- ▶ Added **nvmlDeviceGetLastBBXFlushTime** for retrieving the timestamp and duration of the latest flush of the BBX object to the inforom storage.
- ▶ Added new field Id **NVML\_FI\_DEV\_PCIE\_COUNT\_CORRECTABLE\_ERRORS** for PCIe correctable errors counter.
- ▶ Added new field Id **NVML\_FI\_DEV\_PCIE\_COUNT\_NAKS\_RECEIVED** for PCIe NAK Receive counter.
- ▶ Added new field Id **NVML\_FI\_DEV\_PCIE\_COUNT\_RECEIVER\_ERROR** for PCIe receiver error counter.
- ▶ Added new field Id **NVML\_FI\_DEV\_PCIE\_COUNT\_BAD\_TLP** for PCIe bad TLP counter.
- ▶ Added new field Id **NVML\_FI\_DEV\_PCIE\_COUNT\_NAKS\_SENT** for NAK Send counter.
- ▶ Added new field Id **NVML\_FI\_DEV\_PCIE\_COUNT\_BAD\_DLLP** for PCIe bad DLLP counter.
- ▶ Added new field Id **NVML\_FI\_DEV\_PCIE\_COUNT\_NON\_FATAL\_ERROR** for PCIe non fatal error counter.
- ▶ Added new field Id **NVML\_FI\_DEV\_PCIE\_COUNT\_FATAL\_ERROR** for PCIe fatal error counter.

- ▶ Added new field Id `NVML_FI_DEV_PCIE_COUNT_UNSUPPORTED_REQ` for PCIe unsupported request counter.
- ▶ Added new field Id `NVML_FI_DEV_PCIE_COUNT_LCRC_ERROR` for PCIe LCRC error counter.
- ▶ Added new field Id `NVML_FI_DEV_PCIE_COUNT_LANE_ERROR` for per lane error counter with scope as PCIe lane number.
- ▶ Added `nvmlDeviceGetPowerLimits` for retrieving current Power Limits.
- ▶ Added `nvmlDeviceGetPowerUsage_v2` to retrieve current power usage.
- ▶ Added `nvmlDeviceGetTotalEnergyConsumption_v2` to get current energy consumption.
- ▶ Added `nvmlDeviceSetPowerManagementLimit_v2` to set the power limit.
- ▶ Added new field IDs, `NVML_FI_GPU_POWER_AVERAGE` and `NVML_FI_GPU_POWER_INSTANT`, to query power usage.
- ▶ Renamed `nvmlDeviceCcuGetStreamState` to `nvmlGpmQueryIfStreamingEnabled` and `nvmlDeviceCcuSetStreamState` to `nvmlGpmSetStreamingEnabled`.
- ▶ Added support to display confidential compute protected memory along with fb and bar1 in `nvidia-smi pmon` and `dmon` commands.
- ▶ Added new field IDs `NVML_FI_DEV_TEMPERATURE_SHUTDOWN_TLIMIT`, `NVML_FI_DEV_TEMPERATURE_SLOWDOWN_TLIMIT`, `NVML_FI_DEV_TEMPERATURE_MEM_MAX_TLIMIT`, and `NVML_FI_DEV_TEMPERATURE_GPU_MAX_TLIMIT` to query temperature thresholds on Ada and later architectures.
- ▶ Introduced `ClockEventReasons` and related APIs which should be used instead of `ClockThrottleReasons`. Deprecated `ClockThrottleReasons`.
- ▶ Added ability to get GPS Temperature Threshold with `nvmlDeviceGetTemperatureThreshold` using the new enum `NVML_TEMPERATURE_THRESHOLD_GPS_CURR`.

### Changes between v525 and v530

The following new functionality is exposed on NVIDIA display drivers version 530 Production or later.

- ▶ Fixed a typo in `nvmlGpuP2PStatus_t`: added a new enum entry for `NVML_P2P_STATUS_CHIPSET_NOT_SUPPORTED` with the same numeric value as the existing erroneous entry ("`NVML_P2P_STATUS_CHIPSET_NOT_SUPPORED`").
- ▶ Added `nvmlDeviceGetVgpuSchedulerLog` to fetch the vGPU software scheduler logs.
- ▶ Added `nvmlDeviceGetVgpuSchedulerState` to fetch the vGPU software scheduler state.
- ▶ Added `nvmlDeviceGetVgpuSchedulerCapabilities` to fetch the vGPU software scheduler capabilities.

## Changes between v520 and v525

The following new functionality is exposed on NVIDIA display drivers version 525 Production or later.

- ▶ Added **nvmlDeviceGetPcieAtomicCaps** to report PCIe atomic capabilities.
- ▶ Added **nvmlDeviceCcuGetStreamState** API to report the counter collection unit stream state.
- ▶ Added **nvmlDeviceCcuSetStreamState** API to set the counter collection unit stream state.
- ▶ Removed support for **NVML\_FI\_DEV\_LINK\_SPEED\_MBPS\_L{0..}** field Ids in Hopper. Replaced with **NVML\_FI\_DEV\_NVLINK\_GET\_SPEED** with scope as link Id.
- ▶ Removed support for **NVML\_FI\_DEV\_NVLINK\_CRC\_FLIT\_ERROR\_COUNT{0..}** field Ids in Hopper. Replaced with **NVML\_FI\_DEV\_NVLINK\_ERROR\_DL\_CRC** with scope as link Id.
- ▶ Removed support for **NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_L{0..}** field Ids in Hopper. Replaced with **NVML\_FI\_DEV\_NVLINK\_ERROR\_DL\_REPLAY** with scope as link Id.
- ▶ Removed support for **NVML\_FI\_DEV\_NVLINK\_RECOVERY\_ERROR\_COUNT\_{0..}** field Ids in Hopper. Replaced with **NVML\_FI\_DEV\_NVLINK\_ERROR\_DL\_RECOVERY** with scope as link Id.
- ▶ Added new field ID **NVML\_FI\_DEV\_NVLINK\_GET\_STATE** to get nvlink state.
- ▶ Added new field ID **NVML\_FI\_DEV\_NVLINK\_GET\_VERSION** to get nvlink version.
- ▶ Added new field ID **NVML\_FI\_DEV\_C2C\_LINK\_COUNT** to get C2C link count.
- ▶ Added new field ID **NVML\_FI\_DEV\_C2C\_LINK\_GET\_STATUS** to get C2C link status.
- ▶ Added new field ID **NVML\_FI\_DEV\_C2C\_LINK\_GET\_MAX\_BW** to get C2C link bandwidth.

## Changes between v515 and v520

The following new functionality is exposed on NVIDIA display drivers version 520 Production or later.

- ▶ Added **nvmlDeviceGetMemClkVfOffset** API to report the MemClk VF offset value.
- ▶ Added **nvmlDeviceSetMemClkVfOffset** API to set the MemClk VF offset value.
- ▶ Added **nvmlDeviceGetMemClkMinMaxVfOffset** API to report the Memory clock min and max VF offset that user can set for a specified GPU.
- ▶ Added **nvmlDeviceGetTargetFanSpeed** API to report the intended target speed of the device's specified fan.
- ▶ Added **nvmlDeviceGetGpcClkMinMaxVfOffset** API to report the Graphics clock min and max VF offset that user can set for a specified GPU.
- ▶ Added **nvmlGpmMetricsGet** to calculate GPM metrics from two GPM samples.
- ▶ Added **nvmlGpmSampleFree** to free allocated GPM sample.

- ▶ Added `nvmlGpmSampleAlloc` to allocate a GPM sample.
- ▶ Added `nvmlGpmSampleGet` to retrieve a GPM snapshot.
- ▶ Added `nvmlGpmQueryDeviceSupport` to query whether a device supports GPM
- ▶ Added `nvmlDeviceGetSupportedPowerModes` API to report the GPU's supported power mode mask.
- ▶ Added `nvmlDeviceGetPowerMode` API to report the GPU's current power mode.
- ▶ Added `nvmlDeviceSetPowerMode` API to set the new power mode.
- ▶ Added `nvmlDeviceGetFanControlPolicy_v2` API to report the control policy for a specified GPU fan.
- ▶ Added `nvmlDeviceSetFanControlPolicy` API to set the control policy for a specified GPU fan.

### Changes between v510 and v515

The following new functionality is exposed on NVIDIA display drivers version 515 Production or later.

- ▶ Added `nvmlDeviceGetDefaultECCMode` API to report the GPU's default ECC Mode.
- ▶ Added `nvmlDeviceGetPcieSpeed` API to report the GPU's PCIe link speed.
- ▶ Added `nvmlDeviceGetDynamicPstatesInfo` API to report the GPU's P-states information.
- ▶ Added `nvmlDeviceSetFanSpeed_v2` API to set the GPU's fan speed.
- ▶ Added `nvmlDeviceSetDefaultFanSpeed_v2` API to set the GPU's default fan speed.
- ▶ Added `nvmlDeviceGetThermalSettings` API to report the GPU's thermal system information.
- ▶ Added `nvmlDeviceGetMinMaxClockOfPState` API to report the min and max clocks of some clock domain for a given PState.
- ▶ Added `nvmlDeviceGetSupportedPerformanceStates` API to get all supported Performance States (P-States) for the GPU.
- ▶ Added `nvmlDeviceGetGpcClkVfOffset` API to report the GPCCLK VF offset value.
- ▶ Added `nvmlDeviceSetGpcClkVfOffset` API to set the GPCCLK VF offset value.
- ▶ Added `nvmlDeviceGetMinMaxFanSpeed` API to report the min and max fan speed that user can set for a specified GPU fan.

### Changes between v495 and v510

The following new functionality is exposed on NVIDIA display drivers version 510 Production or later.

- ▶ Added `nvmlDeviceGetGpuInstanceProfileInfoV` and `nvmlGpuInstanceGetComputeInstanceProfileInfoV` APIs to include the profile name in their output.

- ▶ Added `nvmlDeviceGetMemoryBusWidth` API to report the GPU's Memory Bus Width.
- ▶ Added `nvmlDeviceGetPcieLinkMaxSpeed` API to report the GPU's PCIe Max Speed.
- ▶ Added `nvmlDeviceGetPowerSource` API to report the GPU's power source as AC or battery.
- ▶ Added `nvmlDeviceGetNumberOfFans` API to report the GPU's number of fans.
- ▶ Added `nvmlDeviceGetNumGpuCores` API to report the GPU's number of cores.
- ▶ Added `nvmlDeviceGetMemoryInfo_v2`. The new version accounts separately for system-reserved memory, and includes it in the used memory amount. The previous version of the API reduced the total memory amount by the amount of system-reserved memory.
- ▶ Added `nvmlDeviceGetAdaptiveClockInfoStatus` API to report the status of adaptive clocking for the GPU.

### Changes between v465 and v470

The following new functionality is exposed on NVIDIA display drivers version 470 Production or later.

- ▶ Added new MIG GPU instance profile `NVML_GPU_INSTANCE_PROFILE_1_SLICE_REV1`.
- ▶ Added `nvmlDeviceGetGpuInstancePossiblePlacements_v2`. The previous version of the API will not support the profiles with possible placements greater than its total capacity, such as `NVML_GPU_INSTANCE_PROFILE_1_SLICE_REV1`.

### Changes between v460 and v465

The following new functionality is exposed on NVIDIA display drivers version 465 Production or later.

- ▶ Added new `NVML_BRAND_*` enumeration values for `NVIDIA`, `NVIDIA_RTX`, `GEFORCE_RTX`, `QUADRO_RTX` and `TITAN_RTX`.
- ▶ Updated `nvmlDeviceGetHandleByUUID` to make it MIG-aware.
- ▶ Updated `nvmlDeviceGetUUID` to return MIG UUIDs in the canonical format, 'MIG-UUID'.
- ▶ Updated `nvmlDeviceGetHandleByUUID` to accept both UUID formats, 'MIG-UUID' and 'MIG-GPU UUID/GID/CID'.
- ▶ The `nvmlDeviceSetAPIRestriction` and `nvmlDeviceGetAPIRestriction` APIs would no longer support the ability to toggle root-only requirement for `nvmlDeviceSetApplicationsClocks` and `nvmlDeviceResetApplicationsClocks`.

### Changes between v450 and v460

The following new functionality is exposed on NVIDIA display drivers version 460 Production or later.



- ▶ Added `nvmlDeviceCreateGpuInstanceWithPlacement` to allow placement specification when creating a new MIG GPU instance.

### Changes between v445 and v450

The following new functionality is exposed on NVIDIA display drivers version 450 Production or later.

- ▶ Updated `nvmlDeviceGetFanSpeed` and `nvmlDeviceGetFanSpeed_v2` for allowing fan speeds greater than 100% to be reported.
- ▶ Added `nvmlDeviceGetCpuAffinityWithinScope` to determine the closest processor(s) within a NUMA node or socket.
- ▶ Added `nvmlDeviceGetMemoryAffinity` to determine the closest NUMA node(s) within a NUMA node or socket.
- ▶ Added support to query and disable MIG mode on Windows.

### Changes between v418 and v445

The following new functionality is exposed on NVIDIA display drivers version 445 Production or later.

- ▶ Added support for the NVIDIA Ampere architecture.
- ▶ Added support for Multi Instance GPU management. Refer to the "Multi Instance GPU Management" section for details.

### Changes between v361 and v418

The following new functionality is exposed on NVIDIA display drivers version 418 Production or later.

- ▶ Added support for the Volta and Turing architectures, bug fixes, performance improvements, and new features.

### Changes between v349 and v361

The following new functionality is exposed on NVIDIA display drivers version 361 Production or later.

- ▶ Added `nvmlDeviceGetBoardPartNumber` to return GPU part numbers
- ▶ Removed support for exclusive thread compute mode (Deprecated in 7.5)
- ▶ Added `NVML_CLOCK_VIDEO` (encoder/decoder) clock type as a supported clock type for `nvmlDeviceGetClockInfo` and `nvmlDeviceGetMaxClockInfo`.

### Changes between v346 and v349

The following new functionality is exposed on NVIDIA display drivers version 349 Production or later.

- ▶ Added `nvmlDeviceGetTopologyCommonAncestor` to find the common path between two devices.
- ▶ Added `nvmlDeviceGetTopologyNearestGpus` to get a set of GPUs given a path level.
- ▶ Added `nvmlSystemGetTopologyGpuSet` to retrieve a set of GPUs with a given CPU affinity.
- ▶ Discontinued Perl bindings support.
- ▶ Updated `nvmlDeviceGetAccountingPids`, `nvmlDeviceGetAccountingBufferSize` and `nvmlDeviceGetAccountingStats` to report accounting information for both active and terminated processes. The execution time field in `nvmlAccountingStats_t` structure is populated only when the process is terminated.

### Changes between v340 and v346

The following new functionality is exposed on NVIDIA display drivers version 346 Production or later.

- ▶ Added `nvmlDeviceGetGraphicsRunningProcesses` to get information about Graphics Processes running on a device.
- ▶ Added `nvmlDeviceGetPcieReplayCounter` to get PCI replay counters.
- ▶ Added `nvmlDeviceGetPcieThroughput` to get PCI utilization information.
- ▶ Discontinued Perl bindings support.

### Changes between NVML v331 and v340

The following new functionality is exposed on NVIDIA display drivers version 340 Production or later.

- ▶ Added `nvmlDeviceGetSamples` to get recent power, utilization and clock samples for the GPU.
- ▶ Added `nvmlDeviceGetTemperatureThreshold` to get temperature thresholds for the GPU.
- ▶ Added `nvmlDeviceGetBrand` to get the brand name of the GPU.
- ▶ Added `nvmlDeviceGetViolationStatus` to get the duration of time during which the device was throttled (lower than requested clocks) due to power or thermal constraints. Violations due to thermal capping is not supported at this time.
- ▶ Added `nvmlDeviceGetEncoderUtilization` to get the GPU video encoder utilization.
- ▶ Added `nvmlDeviceGetDecoderUtilization` to get the GPU video decoder utilization.
- ▶ Added `nvmlDeviceGetCpuAffinity` to get the closest processor(s) affinity to a particular GPU.
- ▶ Added `nvmlDeviceSetCpuAffinity` to set the affinity of a particular GPU to the closest processor.
- ▶ Added `nvmlDeviceClearCpuAffinity` to clear the affinity of a particular GPU.
- ▶ Added `nvmlDeviceGetBoardId` to get a unique boardId for the running system.

- ▶ Added `nvmlDeviceGetMultiGpuBoard` to get whether the device is on a multiGPU board.
- ▶ Added `nvmlDeviceGetAutoBoostedClocksEnabled` and `nvmlDeviceSetAutoBoostedClocksEnabled` for querying and setting the state of auto boosted clocks on supporting hardware.
- ▶ Added `nvmlDeviceSetDefaultAutoBoostedClocksEnabled` for setting the default state of auto boosted clocks on supporting hardware.

### Changes between NVML v5.319 Update and v331

The following new functionality is exposed on NVIDIA display drivers version 331 or later.

- ▶ Added `nvmlDeviceGetMinorNumber` to get the minor number for the device.
- ▶ Added `nvmlDeviceGetBAR1MemoryInfo` to get BAR1 total, available and used memory size.
- ▶ Added `nvmlDeviceGetBridgeChipInfo` to get the information related to bridge chip firmware.
- ▶ Added enforced power limit query API `nvmlDeviceGetEnforcedPowerLimit`
- ▶ Updated `nvmlEventSetWait` to return xid event data in case of xid error event.

### Changes between NVML v5.319 RC and v5.319 Update

The following new functionality is exposed on NVIDIA display drivers version 319 Update or later.

- ▶ Added `nvmlDeviceSetAPIRestriction` and `nvmlDeviceGetAPIRestriction`, with initial ability to toggle root-only requirement for `nvmlDeviceSetApplicationsClocks` and `nvmlDeviceResetApplicationsClocks`.

### Changes between NVML v4.304 Production and v5.319 RC

The following new functionality is exposed on NVIDIA display drivers version 319 RC or later.

- ▶ Added `_v2` versions of `nvmlDeviceGetHandleByIndex` and `nvmlDeviceGetCount` that also count devices not accessible by current user
  - ▶ `nvmlDeviceGetHandleByIndex_v2` (default) can also return `NVML_ERROR_NO_PERMISSION`
- ▶ Added `nvmlInit_v2` and `nvmlDeviceGetHandleByIndex_v2` that is safer and thus recommended function for initializing the library
  - ▶ `nvmlInit_v2` lazily initializes only requested devices (queried with `nvmlDeviceGetHandle*`)
  - ▶ `nvml.h` defines `nvmlInit_v2` and `nvmlDeviceGetHandleByIndex_v2` as default functions

- ▶ Added `nvmlDeviceGetIndex`
- ▶ Added `NVML_ERROR_GPU_IS_LOST` to report GPUs that have fallen off the bus.
  - ▶ All NVML device APIs can return this error code, as a GPU can fall off the bus at any time.
- ▶ Added new class of APIs for gathering process statistics (`nvmlAccountingStats`)
- ▶ Application Clocks are no longer supported on GPU's from Quadro product line
- ▶ Added APIs to support dynamic page retirement. See `nvmlDeviceGetRetiredPages` and `nvmlDeviceGetRetiredPagesPendingStatus`
- ▶ Renamed `nvmlClocksThrottleReasonUserDefinedClocks` to `nvmlClocksThrottleReasonApplicationsClocksSetting`. Old name is deprecated and can be removed in one of the next major releases.
- ▶ Added `nvmlDeviceGetDisplayActive` and updated documentation to clarify how it differs from `nvmlDeviceGetDisplayMode`

### Changes between NVML v4.304 RC and v4.304 Production

The following new functionality is exposed on NVIDIA display drivers version 304 Production or later.

- ▶ Added `nvmlDeviceGetGpuOperationMode` and `nvmlDeviceSetGpuOperationMode`.

### Changes between NVML v3.295 and v4.304 RC

The following new functionality is exposed on NVIDIA display drivers version 304 RC or later.

- ▶ Added `nvmlDeviceGetInforomConfigurationChecksum` and `nvmlDeviceValidateInforom`.
- ▶ Added `nvmlDeviceGetDisplayActive` and updated documentation to clarify how it differs from `nvmlDeviceGetDisplayMode`.
- ▶ Added new error return value for initialization failure due to kernel module not receiving interrupts.
- ▶ Added `nvmlDeviceSetApplicationsClocks`, `nvmlDeviceGetApplicationsClock`, `nvmlDeviceResetApplicationsClocks`.
- ▶ Added `nvmlDeviceGetSupportedMemoryClocks` and `nvmlDeviceGetSupportedGraphicsClocks`.
- ▶ Added `nvmlDeviceGetPowerManagementLimitConstraints`, `nvmlDeviceGetPowerManagementDefaultLimit` and `nvmlDeviceSetPowerManagementLimit`.
- ▶ Added `nvmlDeviceGetInforomImageVersion`.
- ▶ Expanded `nvmlDeviceGetUUID` to support all CUDA capable GPUs.
- ▶ Deprecated `nvmlDeviceGetDetailedEccErrors` in favor of `nvmlDeviceGetMemoryErrorCounter`.

- ▶ Added NVML\_MEMORY\_LOCATION\_TEXTURE\_MEMORY to support reporting of texture memory error counters.
- ▶ Added `nvmlDeviceGetCurrentClocksThrottleReasons` and `nvmlDeviceGetSupportedClocksThrottleReasons`.
- ▶ NVML\_CLOCK\_SM is now also reported on supported Kepler devices.
- ▶ Dropped support for GT200 based Tesla brand GPUs: C1060, M1060, S1070.

### Changes between NVML v2.285 and v3.295

The following new functionality is exposed on NVIDIA display drivers version 295 or later.

- ▶ Deprecated `nvmlDeviceGetHandleBySerial` in favor of newly added `nvmlDeviceGetHandleByUUID`.
- ▶ Marked the input parameters of `nvmlDeviceGetHandleBySerial`, `nvmlDeviceGetHandleByUUID` and `nvmlDeviceGetHandleByPciBusId` as const.
- ▶ Added `nvmlDeviceOnSameBoard`.
- ▶ Added `nvmlConstants` defines.
- ▶ Added `nvmlDeviceGetMaxPcieLinkGeneration`, `nvmlDeviceGetMaxPcieLinkWidth`, `nvmlDeviceGetCurrPcieLinkGeneration`, `nvmlDeviceGetCurrPcieLinkWidth`.
- ▶ Format change of `nvmlDeviceGetUUID` output to match the UUID standard. This function will return a different value.
- ▶ `nvmlDeviceGetDetailedEccErrors` will report zero for unsupported ECC error counters when a subset of ECC error counters are supported.

### Changes between NVML v1.0 and v2.285

The following new functionality is exposed on NVIDIA display drivers version 285 or later.

- ▶ Added possibility to query separately current and pending driver model with `nvmlDeviceGetDriverModel`.
- ▶ Added API `nvmlDeviceGetVbiosVersion` function to report VBIOS version.
- ▶ Added `pciSubSystemId` to `nvmlPciInfo_t` struct.
- ▶ Added API `nvmlErrorString` function to convert error code to string.
- ▶ Updated docs to indicate we support M2075 and C2075.
- ▶ Added API `nvmlSystemGetHicVersion` function to report HIC firmware version.
- ▶ Added NVML versioning support
  - ▶ Functions that changed API and/or size of structs have appended versioning suffix (e.g., `nvmlDeviceGetPciInfo_v2`). Appropriate C defines have been added that map old function names to the newer version of the function.
- ▶ Added support for concurrent library usage by multiple libraries.

- ▶ Added API `nvmlDeviceGetMaxClockInfo` function for reporting device's clock limits.
- ▶ Added new error code `NVML_ERROR_DRIVER_NOT_LOADED` used by `nvmlInit`.
- ▶ Extended `nvmlPciInfo_t` struct with new field: sub system id.
- ▶ Added NVML support on Windows guest account.
- ▶ Changed format of `pciBusId` string (to `XXXX:XX:XX.X`) of `nvmlPciInfo_t`.
- ▶ Parsing of `busId` in `nvmlDeviceGetHandleByPciBusId` is less restrictive. You can pass `0:2:0.0` or `0000:02:00` and other variations.
- ▶ Added API for events waiting for GPU events (Linux only) see docs of `nvmlEvents`.
- ▶ Added API `nvmlDeviceGetComputeRunningProcesses` and `nvmlSystemGetProcessName` functions for looking up currently running compute applications.
- ▶ Deprecated `nvmlDeviceGetPowerState` in favor of `nvmlDeviceGetPerformanceState`.

# Chapter 4.

## MODULES

Here is a list of all modules:

- ▶ Device Structs
- ▶ Device Enums
- ▶ Field Value Enums
- ▶ Unit Structs
- ▶ Accounting Statistics
- ▶ Encoder Structs
- ▶ Frame Buffer Capture Structures
- ▶ definitions related to the drain state
- ▶ /nvmDevice definitions related to Confidential Computing
- ▶ Initialization and Cleanup
- ▶ Error reporting
- ▶ Constants
- ▶ System Queries
- ▶ Unit Queries
- ▶ Device Queries
  - ▶ CPU and Memory Affinity
- ▶ Unit Commands
- ▶ Device Commands
- ▶ NvLink Methods
- ▶ Event Handling Methods
  - ▶ Event Types
- ▶ Drain states
- ▶ Field Value Queries
- ▶ Enums, Constants and Structs
- ▶ vGPU APIs
- ▶ vGPU Management

- ▶ vGPU Migration
- ▶ vGPU Utilization and Accounting
- ▶ Excluded GPU Queries
- ▶ Multi Instance GPU Management
- ▶ NVML GPM
  - ▶ GPM Enums
  - ▶ GPM Structs
  - ▶ GPM Functions
- ▶ VirtualGPU
  - ▶ vGPU Enums
  - ▶ vGPU Constants
  - ▶ vGPU Structs
- ▶ NvmlClocksEventReasons

## 4.1. Device Structs



```
struct nvmlPciInfoExt_v1_t
struct nvmlPciInfo_t
struct nvmlEccErrorCounts_t
struct nvmlUtilization_t
struct nvmlMemory_t
struct nvmlMemory_v2_t
struct nvmlBAR1Memory_t
struct nvmlProcessInfo_v1_t
struct nvmlProcessInfo_t
struct nvmlProcessDetail_v1_t
struct nvmlProcessDetailList_v1_t
struct nvmlC2cModelInfo_v1_t
struct nvmlRowRemapperHistogramValues_t
struct nvmlNvLinkUtilizationControl_t
struct nvmlBridgeChipInfo_t
struct nvmlBridgeChipHierarchy_t
union nvmlValue_t
struct nvmlSample_t
```

## struct nvmlViolationTime\_t

## enum nvmlBridgeChipType\_t

Enum to represent type of bridge chip

### Values

NVML\_BRIDGE\_CHIP\_PLX = 0  
 NVML\_BRIDGE\_CHIP\_BRO4 = 1

## enum nvmlNvLinkUtilizationCountUnits\_t

Enum to represent the NvLink utilization counter packet units

### Values

NVML\_NVLINK\_COUNTER\_UNIT\_CYCLES = 0  
 NVML\_NVLINK\_COUNTER\_UNIT\_PACKETS = 1  
 NVML\_NVLINK\_COUNTER\_UNIT\_BYTES = 2  
 NVML\_NVLINK\_COUNTER\_UNIT\_RESERVED = 3  
 NVML\_NVLINK\_COUNTER\_UNIT\_COUNT

## enum nvmlNvLinkUtilizationCountPktTypes\_t

Enum to represent the NvLink utilization counter packet types to count \*\* this is ONLY applicable with the units as packets or bytes \*\* as specified in nvmlNvLinkUtilizationCountUnits\_t \*\* all packet filter descriptions are target GPU centric \*\* these can be "OR'd" together

### Values

NVML\_NVLINK\_COUNTER\_PKTFILTER\_NOP = 0x1  
 NVML\_NVLINK\_COUNTER\_PKTFILTER\_READ = 0x2  
 NVML\_NVLINK\_COUNTER\_PKTFILTER\_WRITE = 0x4  
 NVML\_NVLINK\_COUNTER\_PKTFILTER\_RATOM = 0x8  
 NVML\_NVLINK\_COUNTER\_PKTFILTER\_NRATOM = 0x10  
 NVML\_NVLINK\_COUNTER\_PKTFILTER\_FLUSH = 0x20  
 NVML\_NVLINK\_COUNTER\_PKTFILTER\_RESPDATA = 0x40  
 NVML\_NVLINK\_COUNTER\_PKTFILTER\_RESPNODATA = 0x80  
 NVML\_NVLINK\_COUNTER\_PKTFILTER\_ALL = 0xFF

## enum nvmlNvLinkCapability\_t

Enum to represent NvLink queryable capabilities

**Values**

```

NVML_NVLINK_CAP_P2P_SUPPORTED = 0
NVML_NVLINK_CAP_SYSTEMEM_ACCESS = 1
NVML_NVLINK_CAP_P2P_ATOMICS = 2
NVML_NVLINK_CAP_SYSTEMEM_ATOMICS = 3
NVML_NVLINK_CAP_SLI_BRIDGE = 4
NVML_NVLINK_CAP_VALID = 5
NVML_NVLINK_CAP_COUNT

```

**enum nvmlNvLinkErrorCounter\_t**

Enum to represent NvLink queryable error counters

**Values**

```

NVML_NVLINK_ERROR_DL_REPLAY = 0
NVML_NVLINK_ERROR_DL_RECOVERY = 1
NVML_NVLINK_ERROR_DL_CRC_FLIT = 2
NVML_NVLINK_ERROR_DL_CRC_DATA = 3
NVML_NVLINK_ERROR_DL_ECC_DATA = 4
NVML_NVLINK_ERROR_COUNT

```

**enum nvmlIntNvLinkDeviceType\_t**

Enum to represent NvLink's remote device type

**Values**

```

NVML_NVLINK_DEVICE_TYPE_GPU = 0x00
NVML_NVLINK_DEVICE_TYPE_IBMNPU = 0x01
NVML_NVLINK_DEVICE_TYPE_SWITCH = 0x02
NVML_NVLINK_DEVICE_TYPE_UNKNOWN = 0xFF

```

**enum nvmlGpuTopologyLevel\_t**

Represents level relationships within a system between two GPUs The enums are spaced to allow for future relationships

**Values**

```

NVML_TOPOLOGY_INTERNAL = 0
NVML_TOPOLOGY_SINGLE = 10
NVML_TOPOLOGY_MULTIPLE = 20
NVML_TOPOLOGY_HOSTBRIDGE = 30
NVML_TOPOLOGY_NODE = 40
NVML_TOPOLOGY_SYSTEM = 50

```

## enum nvmlSamplingType\_t

Represents Type of Sampling Event

### Values

**NVML\_TOTAL\_POWER\_SAMPLES = 0**

To represent total power drawn by GPU.

**NVML\_GPU\_UTILIZATION\_SAMPLES = 1**

To represent percent of time during which one or more kernels was executing on the GPU.

**NVML\_MEMORY\_UTILIZATION\_SAMPLES = 2**

To represent percent of time during which global (device) memory was being read or written.

**NVML\_ENC\_UTILIZATION\_SAMPLES = 3**

To represent percent of time during which NVENC remains busy.

**NVML\_DEC\_UTILIZATION\_SAMPLES = 4**

To represent percent of time during which NVDEC remains busy.

**NVML\_PROCESSOR\_CLK\_SAMPLES = 5**

To represent processor clock samples.

**NVML\_MEMORY\_CLK\_SAMPLES = 6**

To represent memory clock samples.

**NVML\_MODULE\_POWER\_SAMPLES = 7**

To represent module power samples for total module starting Grace Hopper.

**NVML\_JPG\_UTILIZATION\_SAMPLES = 8**

To represent percent of time during which NVJPG remains busy.

**NVML\_OFA\_UTILIZATION\_SAMPLES = 9**

To represent percent of time during which NVOFA remains busy.

**NVML\_SAMPLINGTYPE\_COUNT**

## enum nvmlPcieUtilCounter\_t

Represents the queryable PCIe utilization counters

### Values

**NVML\_PCIE\_UTIL\_TX\_BYTES = 0**

**NVML\_PCIE\_UTIL\_RX\_BYTES = 1**

**NVML\_PCIE\_UTIL\_COUNT**

## enum nvmlValueType\_t

Represents the type for sample value returned

**Values**

```

NVML_VALUE_TYPE_DOUBLE = 0
NVML_VALUE_TYPE_UNSIGNED_INT = 1
NVML_VALUE_TYPE_UNSIGNED_LONG = 2
NVML_VALUE_TYPE_UNSIGNED_LONG_LONG = 3
NVML_VALUE_TYPE_SIGNED_LONG_LONG = 4
NVML_VALUE_TYPE_SIGNED_INT = 5
NVML_VALUE_TYPE_COUNT

```

**enum nvmlPerfPolicyType\_t**

Represents type of perf policy for which violation times can be queried

**Values**

```

NVML_PERF_POLICY_POWER = 0
    How long did power violations cause the GPU to be below application clocks.
NVML_PERF_POLICY_THERMAL = 1
    How long did thermal violations cause the GPU to be below application clocks.
NVML_PERF_POLICY_SYNC_BOOST = 2
    How long did sync boost cause the GPU to be below application clocks.
NVML_PERF_POLICY_BOARD_LIMIT = 3
    How long did the board limit cause the GPU to be below application clocks.
NVML_PERF_POLICY_LOW_UTILIZATION = 4
    How long did low utilization cause the GPU to be below application clocks.
NVML_PERF_POLICY_RELIABILITY = 5
    How long did the board reliability limit cause the GPU to be below application
    clocks.
NVML_PERF_POLICY_TOTAL_APP_CLOCKS = 10
    Total time the GPU was held below application clocks by any limiter (0 - 5 above).
NVML_PERF_POLICY_TOTAL_BASE_CLOCKS = 11
    Total time the GPU was held below base clocks.
NVML_PERF_POLICY_COUNT

```

**#define NVML\_VALUE\_NOT\_AVAILABLE (-1)**

Special constant that some fields take when they are not available. Used when only part of the struct is not available.

Each structure explicitly states when to check for this value.

**#define NVML\_DEVICE\_PCI\_BUS\_ID\_BUFFER\_SIZE 32**

Buffer size guaranteed to be large enough for pci bus id

```
#define NVML_DEVICE_PCI_BUS_ID_BUFFER_V2_SIZE 16
```

Buffer size guaranteed to be large enough for pci bus id for busIdLegacy

```
#define NVML_DEVICE_PCI_BUS_ID_LEGACY_FMT "%04X:
%02X:%02X.0"
```

PCI format string for busIdLegacy

```
#define NVML_DEVICE_PCI_BUS_ID_FMT "%08X:%02X:
%02X.0"
```

PCI format string for busId

```
#define NVML_DEVICE_PCI_BUS_ID_FMT_ARGS (pciInfo)-
>domain, \ (pciInfo)->bus, \ (pciInfo)->device
```

Utility macro for filling the pci bus id format from a `nvmlPciInfo_t`

```
#define nvmlProcessDetailList_v1
NVML_STRUCT_VERSION(ProcessDetailList, 1)
```

`nvmlProcessDetailList` version

```
#define NVML_NVLINK_MAX_LINKS 18
```

Maximum number of NvLink links supported

```
#define NVML_MAX_PHYSICAL_BRIDGE (128)
```

Maximum limit on Physical Bridges per Board

## 4.2. Device Enums

`struct nvmlClkMonFaultInfo_t`

`struct nvmlClkMonStatus_t`

`enum nvmlEnableState_t`

Generic enable/disable enum.

#### Values

`NVML_FEATURE_DISABLED = 0`

Feature disabled.

`NVML_FEATURE_ENABLED = 1`

Feature enabled.

`enum nvmlBrandType_t`

\* The Brand of the GPU

#### Values

`NVML_BRAND_UNKNOWN = 0`

`NVML_BRAND_QUADRO = 1`

`NVML_BRAND_TESLA = 2`

`NVML_BRAND_NVS = 3`

`NVML_BRAND_GRID = 4`

`NVML_BRAND_GEFORCE = 5`

`NVML_BRAND_TITAN = 6`

`NVML_BRAND_NVIDIA_VAPPS = 7`

`NVML_BRAND_NVIDIA_VPC = 8`

`NVML_BRAND_NVIDIA_VCS = 9`

`NVML_BRAND_NVIDIA_VWS = 10`

`NVML_BRAND_NVIDIA_CLOUD_GAMING = 11`

`NVML_BRAND_NVIDIA_VGAMING =`

`NVML_BRAND_NVIDIA_CLOUD_GAMING`

`NVML_BRAND_QUADRO_RTX = 12`

`NVML_BRAND_NVIDIA_RTX = 13`

`NVML_BRAND_NVIDIA = 14`

`NVML_BRAND_GEFORCE_RTX = 15`

`NVML_BRAND_TITAN_RTX = 16`

`NVML_BRAND_COUNT`

## enum nvmlTemperatureThresholds\_t

Temperature thresholds.

### Values

**NVML\_TEMPERATURE\_THRESHOLD\_SHUTDOWN = 0**  
**NVML\_TEMPERATURE\_THRESHOLD\_SLOWDOWN = 1**  
**NVML\_TEMPERATURE\_THRESHOLD\_MEM\_MAX = 2**  
**NVML\_TEMPERATURE\_THRESHOLD\_GPU\_MAX = 3**  
**NVML\_TEMPERATURE\_THRESHOLD\_ACOUSTIC\_MIN = 4**  
**NVML\_TEMPERATURE\_THRESHOLD\_ACOUSTIC\_CURR = 5**  
**NVML\_TEMPERATURE\_THRESHOLD\_ACOUSTIC\_MAX = 6**  
**NVML\_TEMPERATURE\_THRESHOLD\_COUNT**

## enum nvmlTemperatureSensors\_t

Temperature sensors.

### Values

**NVML\_TEMPERATURE\_GPU = 0**  
 Temperature sensor for the GPU die.  
**NVML\_TEMPERATURE\_COUNT**

## enum nvmlComputeMode\_t

Compute mode.

**NVML\_COMPUTEMODE\_EXCLUSIVE\_PROCESS** was added in CUDA 4.0.  
 Earlier CUDA versions supported a single exclusive mode, which is equivalent to **NVML\_COMPUTEMODE\_EXCLUSIVE\_THREAD** in CUDA 4.0 and beyond.

### Values

**NVML\_COMPUTEMODE\_DEFAULT = 0**  
 Default compute mode -- multiple contexts per device.  
**NVML\_COMPUTEMODE\_EXCLUSIVE\_THREAD = 1**  
 Support Removed.  
**NVML\_COMPUTEMODE\_PROHIBITED = 2**  
 Compute-prohibited mode -- no contexts per device.  
**NVML\_COMPUTEMODE\_EXCLUSIVE\_PROCESS = 3**  
 Compute-exclusive-process mode -- only one context per device, usable from multiple threads at a time.  
**NVML\_COMPUTEMODE\_COUNT**



## enum nvmlMemoryErrorType\_t

Memory error types

### Values

**NVML\_MEMORY\_ERROR\_TYPE\_CORRECTED = 0**

A memory error that was corrected. For ECC errors, these are single bit errors. For Texture memory, these are errors fixed by resend.

**NVML\_MEMORY\_ERROR\_TYPE\_UNCORRECTED = 1**

A memory error that was not corrected. For ECC errors, these are double bit errors. For Texture memory, these are errors where the resend fails.

**NVML\_MEMORY\_ERROR\_TYPE\_COUNT**

Count of memory error types.

## enum nvmlEccCounterType\_t

ECC counter types.

Note: Volatile counts are reset each time the driver loads. On Windows this is once per boot. On Linux this can be more frequent. On Linux the driver unloads when no active clients exist. If persistence mode is enabled or there is always a driver client active (e.g. X11), then Linux also sees per-boot behavior. If not, volatile counts are reset each time a compute app is run.

### Values

**NVML\_VOLATILE\_ECC = 0**

Volatile counts are reset each time the driver loads.

**NVML\_AGGREGATE\_ECC = 1**

Aggregate counts persist across reboots (i.e. for the lifetime of the device).

**NVML\_ECC\_COUNTER\_TYPE\_COUNT**

Count of memory counter types.

## enum nvmlClockType\_t

Clock types.

All speeds are in Mhz.

### Values

**NVML\_CLOCK\_GRAPHICS = 0**

Graphics clock domain.

**NVML\_CLOCK\_SM = 1**

SM clock domain.

**NVML\_CLOCK\_MEM = 2**

Memory clock domain.

**NVML\_CLOCK\_VIDEO = 3**

Video encoder/decoder clock domain.

**NVML\_CLOCK\_COUNT**

Count of clock types.

## enum nvmlClockId\_t

Clock Ids. These are used in combination with `nvmlClockType_t` to specify a single clock value.

### Values

**NVML\_CLOCK\_ID\_CURRENT = 0**

Current actual clock value.

**NVML\_CLOCK\_ID\_APP\_CLOCK\_TARGET = 1**

Target application clock.

**NVML\_CLOCK\_ID\_APP\_CLOCK\_DEFAULT = 2**

Default application clock target.

**NVML\_CLOCK\_ID\_CUSTOMER\_BOOST\_MAX = 3**

OEM-defined maximum clock rate.

**NVML\_CLOCK\_ID\_COUNT**

Count of Clock Ids.

## enum nvmlDriverModel\_t

Driver models.

Windows only.

### Values

**NVML\_DRIVER\_WDDM = 0**

WDDM driver model -- GPU treated as a display device.

**NVML\_DRIVER\_WDM = 1**

WDM (TCC) model (recommended) -- GPU treated as a generic device.

## enum nvmlPstates\_t

Allowed PStates.

### Values

**NVML\_PSTATE\_0 = 0**

Performance state 0 -- Maximum Performance.

**NVML\_PSTATE\_1 = 1**  
Performance state 1.

**NVML\_PSTATE\_2 = 2**  
Performance state 2.

**NVML\_PSTATE\_3 = 3**  
Performance state 3.

**NVML\_PSTATE\_4 = 4**  
Performance state 4.

**NVML\_PSTATE\_5 = 5**  
Performance state 5.

**NVML\_PSTATE\_6 = 6**  
Performance state 6.

**NVML\_PSTATE\_7 = 7**  
Performance state 7.

**NVML\_PSTATE\_8 = 8**  
Performance state 8.

**NVML\_PSTATE\_9 = 9**  
Performance state 9.

**NVML\_PSTATE\_10 = 10**  
Performance state 10.

**NVML\_PSTATE\_11 = 11**  
Performance state 11.

**NVML\_PSTATE\_12 = 12**  
Performance state 12.

**NVML\_PSTATE\_13 = 13**  
Performance state 13.

**NVML\_PSTATE\_14 = 14**  
Performance state 14.

**NVML\_PSTATE\_15 = 15**  
Performance state 15 -- Minimum Performance.

**NVML\_PSTATE\_UNKNOWN = 32**  
Unknown performance state.

## enum nvmlGpuOperationMode\_t

GPU Operation Mode

GOM allows to reduce power usage and optimize GPU throughput by disabling GPU features.

Each GOM is designed to meet specific user needs.

**Values****NVML\_GOM\_ALL\_ON = 0**

Everything is enabled and running at full speed.

**NVML\_GOM\_COMPUTE = 1**

Designed for running only compute tasks. Graphics operations are not allowed

**NVML\_GOM\_LOW\_DP = 2**

Designed for running graphics applications that don't require high bandwidth double precision

**enum nvmlInforomObject\_t**

Available infoROM objects.

**Values****NVML\_INFOTROM\_OEM = 0**

An object defined by OEM.

**NVML\_INFOTROM\_ECC = 1**

The ECC object determining the level of ECC support.

**NVML\_INFOTROM\_POWER = 2**

The power management object.

**NVML\_INFOTROM\_COUNT**

This counts the number of infoROM objects the driver knows about.

**enum nvmlReturn\_t**

Return values for NVML API calls.

**Values****NVML\_SUCCESS = 0**

The operation was successful.

**NVML\_ERROR\_UNINITIALIZED = 1**

NVML was not first initialized with nvmlInit().

**NVML\_ERROR\_INVALID\_ARGUMENT = 2**

A supplied argument is invalid.

**NVML\_ERROR\_NOT\_SUPPORTED = 3**

The requested operation is not available on target device.

**NVML\_ERROR\_NO\_PERMISSION = 4**

The current user does not have permission for operation.

**NVML\_ERROR\_ALREADY\_INITIALIZED = 5**

Deprecated: Multiple initializations are now allowed through ref counting.

**NVML\_ERROR\_NOT\_FOUND = 6**

A query to find an object was unsuccessful.

**NVML\_ERROR\_INSUFFICIENT\_SIZE = 7**

An input argument is not large enough.

**NVML\_ERROR\_INSUFFICIENT\_POWER = 8**

A device's external power cables are not properly attached.

**NVML\_ERROR\_DRIVER\_NOT\_LOADED = 9**

NVIDIA driver is not loaded.

**NVML\_ERROR\_TIMEOUT = 10**

User provided timeout passed.

**NVML\_ERROR\_IRQ\_ISSUE = 11**

NVIDIA Kernel detected an interrupt issue with a GPU.

**NVML\_ERROR\_LIBRARY\_NOT\_FOUND = 12**

NVML Shared Library couldn't be found or loaded.

**NVML\_ERROR\_FUNCTION\_NOT\_FOUND = 13**

Local version of NVML doesn't implement this function.

**NVML\_ERROR\_CORRUPTED\_INFOROM = 14**

infoROM is corrupted

**NVML\_ERROR\_GPU\_IS\_LOST = 15**

The GPU has fallen off the bus or has otherwise become inaccessible.

**NVML\_ERROR\_RESET\_REQUIRED = 16**

The GPU requires a reset before it can be used again.

**NVML\_ERROR\_OPERATING\_SYSTEM = 17**

The GPU control device has been blocked by the operating system/cgroups.

**NVML\_ERROR\_LIB\_RM\_VERSION\_MISMATCH = 18**

RM detects a driver/library version mismatch.

**NVML\_ERROR\_IN\_USE = 19**

An operation cannot be performed because the GPU is currently in use.

**NVML\_ERROR\_MEMORY = 20**

Insufficient memory.

**NVML\_ERROR\_NO\_DATA = 21**

No data.

**NVML\_ERROR\_VGPU\_ECC\_NOT\_SUPPORTED = 22**

The requested vgpu operation is not available on target device, because ECC is enabled.

**NVML\_ERROR\_INSUFFICIENT\_RESOURCES = 23**

Ran out of critical resources, other than memory.

**NVML\_ERROR\_FREQ\_NOT\_SUPPORTED = 24**

Ran out of critical resources, other than memory.

**NVML\_ERROR\_ARGUMENT\_VERSION\_MISMATCH = 25**

The provided version is invalid/unsupported.

**NVML\_ERROR\_DEPRECATED = 26**

The requested functionality has been deprecated.

**NVML\_ERROR\_NOT\_READY = 27**

The system is not ready for the request.

**NVML\_ERROR\_GPU\_NOT\_FOUND = 28**

No GPUs were found.

**NVML\_ERROR\_INVALID\_STATE = 29**

Resource not in correct state to perform requested operation.

**NVML\_ERROR\_UNKNOWN = 999**

An internal driver error occurred.

## enum nvmlMemoryLocation\_t

See [nvmlDeviceGetMemoryErrorCounter](#)

### Values

**NVML\_MEMORY\_LOCATION\_L1\_CACHE = 0**

GPU L1 Cache.

**NVML\_MEMORY\_LOCATION\_L2\_CACHE = 1**

GPU L2 Cache.

**NVML\_MEMORY\_LOCATION\_DRAM = 2**

Turing+ DRAM.

**NVML\_MEMORY\_LOCATION\_DEVICE\_MEMORY = 2**

GPU Device Memory.

**NVML\_MEMORY\_LOCATION\_REGISTER\_FILE = 3**

GPU Register File.

**NVML\_MEMORY\_LOCATION\_TEXTURE\_MEMORY = 4**

GPU Texture Memory.

**NVML\_MEMORY\_LOCATION\_TEXTURE\_SHM = 5**

Shared memory.

**NVML\_MEMORY\_LOCATION\_CBU = 6**

CBU.

**NVML\_MEMORY\_LOCATION\_SRAM = 7**

Turing+ SRAM.

**NVML\_MEMORY\_LOCATION\_COUNT**

This counts the number of memory locations the driver knows about.

## enum nvmlPageRetirementCause\_t

Causes for page retirement

### Values

**NVML\_PAGE\_RETIREMENT\_CAUSE\_MULTIPLE\_SINGLE\_BIT\_ECC\_ERRORS = 0**

Page was retired due to multiple single bit ECC error.

**NVML\_PAGE\_RETIREMENT\_CAUSE\_DOUBLE\_BIT\_ECC\_ERROR = 1**

Page was retired due to double bit ECC error.

**NVML\_PAGE\_RETIREMENT\_CAUSE\_COUNT**

## enum nvmlRestrictedAPI\_t

API types that allow changes to default permission restrictions

### Values

**NVML\_RESTRICTED\_API\_SET\_APPLICATION\_CLOCKS = 0**

APIs that change application clocks, see `nvmlDeviceSetApplicationsClocks` and see `nvmlDeviceResetApplicationsClocks`

**NVML\_RESTRICTED\_API\_SET\_AUTO\_BOOSTED\_CLOCKS = 1**

APIs that enable/disable Auto Boosted clocks see `nvmlDeviceSetAutoBoostedClocksEnabled`

**NVML\_RESTRICTED\_API\_COUNT**

## #define nvmlFlagDefault 0x00

Generic flag used to specify the default behavior of some functions. See description of particular functions for details.

## #define nvmlFlagForce 0x01

Generic flag used to force some behavior. See description of particular functions for details.

## #define MAX\_CLK\_DOMAINS 32

Max Clock Monitors available

## #define nvmlEccBitType\_t nvmlMemoryErrorType\_t

ECC bit types.

Deprecated See `nvmlMemoryErrorType_t` for a more flexible type

## #define NVML\_SINGLE\_BIT\_ECC

**NVML\_MEMORY\_ERROR\_TYPE\_CORRECTED**

Single bit ECC errors

Deprecated Mapped to `NVML_MEMORY_ERROR_TYPE_CORRECTED`

## #define NVML\_DOUBLE\_BIT\_ECC

**NVML\_MEMORY\_ERROR\_TYPE\_UNCORRECTED**

Double bit ECC errors

Deprecated Mapped to NVML\_MEMORY\_ERROR\_TYPE\_UNCORRECTED

## 4.3. Field Value Enums

```
struct nvmlFieldValue_t
```

```
#define NVML_FI_DEV_ECC_CURRENT 1
```

Current ECC mode. 1=Active. 0=Inactive.

Field Identifiers.

All Identifiers pertain to a device. Each ID is only used once and is guaranteed never to change.

```
#define NVML_FI_DEV_ECC_PENDING 2
```

Pending ECC mode. 1=Active. 0=Inactive.

```
#define NVML_FI_DEV_ECC_SBE_VOL_TOTAL 3
```

Total single bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_DBE_VOL_TOTAL 4
```

Total double bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_SBE_AGG_TOTAL 5
```

Total single bit aggregate (persistent) ECC errors.

```
#define NVML_FI_DEV_ECC_DBE_AGG_TOTAL 6
```

Total double bit aggregate (persistent) ECC errors.

```
#define NVML_FI_DEV_ECC_SBE_VOL_L1 7
```

L1 cache single bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_DBE_VOL_L1 8
```

L1 cache double bit volatile ECC errors.



```
#define NVML_FI_DEV_ECC_SBE_VOL_L2 9
```

L2 cache single bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_DBE_VOL_L2 10
```

L2 cache double bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_SBE_VOL_DEV 11
```

Device memory single bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_DBE_VOL_DEV 12
```

Device memory double bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_SBE_VOL_REG 13
```

Register file single bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_DBE_VOL_REG 14
```

Register file double bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_SBE_VOL_TEX 15
```

Texture memory single bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_DBE_VOL_TEX 16
```

Texture memory double bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_DBE_VOL_CBU 17
```

CBU double bit volatile ECC errors.

```
#define NVML_FI_DEV_ECC_SBE_AGG_L1 18
```

L1 cache single bit aggregate (persistent) ECC errors.

```
#define NVML_FI_DEV_ECC_DBE_AGG_L1 19
```

L1 cache double bit aggregate (persistent) ECC errors.

**#define NVML\_FI\_DEV\_ECC\_SBE\_AGG\_L2 20**

L2 cache single bit aggregate (persistent) ECC errors.

**#define NVML\_FI\_DEV\_ECC\_DBE\_AGG\_L2 21**

L2 cache double bit aggregate (persistent) ECC errors.

**#define NVML\_FI\_DEV\_ECC\_SBE\_AGG\_DEV 22**

Device memory single bit aggregate (persistent) ECC errors.

**#define NVML\_FI\_DEV\_ECC\_DBE\_AGG\_DEV 23**

Device memory double bit aggregate (persistent) ECC errors.

**#define NVML\_FI\_DEV\_ECC\_SBE\_AGG\_REG 24**

Register File single bit aggregate (persistent) ECC errors.

**#define NVML\_FI\_DEV\_ECC\_DBE\_AGG\_REG 25**

Register File double bit aggregate (persistent) ECC errors.

**#define NVML\_FI\_DEV\_ECC\_SBE\_AGG\_TEX 26**

Texture memory single bit aggregate (persistent) ECC errors.

**#define NVML\_FI\_DEV\_ECC\_DBE\_AGG\_TEX 27**

Texture memory double bit aggregate (persistent) ECC errors.

**#define NVML\_FI\_DEV\_ECC\_DBE\_AGG\_CBU 28**

CBU double bit aggregate ECC errors.

**#define NVML\_FI\_DEV\_RETIRED\_SBE 29**

Number of retired pages because of single bit errors.

**#define NVML\_FI\_DEV\_RETIRED\_DBE 30**

Number of retired pages because of double bit errors.

```
#define NVML_FI_DEV_RETIRED_PENDING 31
```

If any pages are pending retirement. 1=yes. 0=no.

```
#define
```

```
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L0 32
```

NVLink flow control CRC Error Counter for Lane 0.

```
#define
```

```
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L1 33
```

NVLink flow control CRC Error Counter for Lane 1.

```
#define
```

```
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L2 34
```

NVLink flow control CRC Error Counter for Lane 2.

```
#define
```

```
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L3 35
```

NVLink flow control CRC Error Counter for Lane 3.

```
#define
```

```
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L4 36
```

NVLink flow control CRC Error Counter for Lane 4.

```
#define
```

```
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_L5 37
```

NVLink flow control CRC Error Counter for Lane 5.

```
#define
```

```
NVML_FI_DEV_NVLINK_CRC_FLIT_ERROR_COUNT_TOTAL
```

```
38
```

NVLink flow control CRC Error Counter total for all Lanes.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L0 39**

NVLink data CRC Error Counter for Lane 0.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L1 40**

NVLink data CRC Error Counter for Lane 1.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L2 41**

NVLink data CRC Error Counter for Lane 2.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L3 42**

NVLink data CRC Error Counter for Lane 3.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L4 43**

NVLink data CRC Error Counter for Lane 4.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L5 44**

NVLink data CRC Error Counter for Lane 5.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_TOTAL  
45**

NvLink data CRC Error Counter total for all Lanes.

**#define**

**NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_L0 46**

NVLink Replay Error Counter for Lane 0.

**#define**

**NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_L1 47**

NVLink Replay Error Counter for Lane 1.

**#define**

**NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_L2 48**

NVLink Replay Error Counter for Lane 2.

**#define**

**NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_L3 49**

NVLink Replay Error Counter for Lane 3.

**#define**

**NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_L4 50**

NVLink Replay Error Counter for Lane 4.

**#define**

**NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_L5 51**

NVLink Replay Error Counter for Lane 5.

**#define**

**NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_TOTAL 52**

NVLink Replay Error Counter total for all Lanes.

**#define**

**NVML\_FI\_DEV\_NVLINK\_RECOVERY\_ERROR\_COUNT\_L0 53**

NVLink Recovery Error Counter for Lane 0.

**#define**

**NVML\_FI\_DEV\_NVLINK\_RECOVERY\_ERROR\_COUNT\_L1 54**

NVLink Recovery Error Counter for Lane 1.

```
#define  
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L2 55  
NVLink Recovery Error Counter for Lane 2.
```

```
#define  
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L3 56  
NVLink Recovery Error Counter for Lane 3.
```

```
#define  
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L4 57  
NVLink Recovery Error Counter for Lane 4.
```

```
#define  
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L5 58  
NVLink Recovery Error Counter for Lane 5.
```

```
#define  
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_TOTAL  
59  
NVLink Recovery Error Counter total for all Lanes.
```

```
#define NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L0 60  
NVLink Bandwidth Counter for Counter Set 0, Lane 0.
```

```
#define NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L1 61  
NVLink Bandwidth Counter for Counter Set 0, Lane 1.
```

```
#define NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L2 62  
NVLink Bandwidth Counter for Counter Set 0, Lane 2.
```

```
#define NVML_FI_DEV_NVLINK_BANDWIDTH_C0_L3 63  
NVLink Bandwidth Counter for Counter Set 0, Lane 3.
```

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C0\_L4 64**

NVLink Bandwidth Counter for Counter Set 0, Lane 4.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C0\_L5 65**

NVLink Bandwidth Counter for Counter Set 0, Lane 5.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C0\_TOTAL  
66**

NVLink Bandwidth Counter Total for Counter Set 0, All Lanes.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C1\_L0 67**

NVLink Bandwidth Counter for Counter Set 1, Lane 0.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C1\_L1 68**

NVLink Bandwidth Counter for Counter Set 1, Lane 1.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C1\_L2 69**

NVLink Bandwidth Counter for Counter Set 1, Lane 2.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C1\_L3 70**

NVLink Bandwidth Counter for Counter Set 1, Lane 3.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C1\_L4 71**

NVLink Bandwidth Counter for Counter Set 1, Lane 4.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C1\_L5 72**

NVLink Bandwidth Counter for Counter Set 1, Lane 5.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C1\_TOTAL  
73**

NVLink Bandwidth Counter Total for Counter Set 1, All Lanes.

**#define NVML\_FI\_DEV\_PERF\_POLICY\_POWER 74**

Perf Policy Counter for Power Policy.

```
#define NVML_FI_DEV_PERF_POLICY_THERMAL 75
```

Perf Policy Counter for Thermal Policy.

```
#define NVML_FI_DEV_PERF_POLICY_SYNC_BOOST 76
```

Perf Policy Counter for Sync boost Policy.

```
#define NVML_FI_DEV_PERF_POLICY_BOARD_LIMIT 77
```

Perf Policy Counter for Board Limit.

```
#define NVML_FI_DEV_PERF_POLICY_LOW_UTILIZATION  
78
```

Perf Policy Counter for Low GPU Utilization Policy.

```
#define NVML_FI_DEV_PERF_POLICY_RELIABILITY 79
```

Perf Policy Counter for Reliability Policy.

```
#define
```

```
NVML_FI_DEV_PERF_POLICY_TOTAL_APP_CLOCKS 80
```

Perf Policy Counter for Total App Clock Policy.

```
#define
```

```
NVML_FI_DEV_PERF_POLICY_TOTAL_BASE_CLOCKS 81
```

Perf Policy Counter for Total Base Clocks Policy.

```
#define NVML_FI_DEV_MEMORY_TEMP 82
```

Memory temperature for the device.

```
#define NVML_FI_DEV_TOTAL_ENERGY_CONSUMPTION 83
```

Total energy consumption for the GPU in mJ since the driver was last reloaded.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L0 84
```

NVLink Speed in MBps for Link 0.



```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L1 85
```

NVLink Speed in MBps for Link 1.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L2 86
```

NVLink Speed in MBps for Link 2.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L3 87
```

NVLink Speed in MBps for Link 3.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L4 88
```

NVLink Speed in MBps for Link 4.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L5 89
```

NVLink Speed in MBps for Link 5.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_COMMON 90
```

Common NVLink Speed in MBps for active links.

```
#define NVML_FI_DEV_NVLINK_LINK_COUNT 91
```

Number of NVLinks present on the device.

```
#define NVML_FI_DEV_RETIRED_PENDING_SBE 92
```

If any pages are pending retirement due to SBE. 1=yes. 0=no.

```
#define NVML_FI_DEV_RETIRED_PENDING_DBE 93
```

If any pages are pending retirement due to DBE. 1=yes. 0=no.

```
#define NVML_FI_DEV_PCIE_REPLAY_COUNTER 94
```

PCIe replay counter.

```
#define
```

```
NVML_FI_DEV_PCIE_REPLAY_ROLLOVER_COUNTER 95
```

PCIe replay rollover counter.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_FLIT\_ERROR\_COUNT\_L6 96**

NVLink flow control CRC Error Counter for Lane 6.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_FLIT\_ERROR\_COUNT\_L7 97**

NVLink flow control CRC Error Counter for Lane 7.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_FLIT\_ERROR\_COUNT\_L8 98**

NVLink flow control CRC Error Counter for Lane 8.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_FLIT\_ERROR\_COUNT\_L9 99**

NVLink flow control CRC Error Counter for Lane 9.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_FLIT\_ERROR\_COUNT\_L10  
100**

NVLink flow control CRC Error Counter for Lane 10.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_FLIT\_ERROR\_COUNT\_L11  
101**

NVLink flow control CRC Error Counter for Lane 11.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L6 102**

NVLink data CRC Error Counter for Lane 6.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L7 103**

NVLink data CRC Error Counter for Lane 7.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L8 104**

NVLink data CRC Error Counter for Lane 8.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L9 105**

NVLink data CRC Error Counter for Lane 9.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L10  
106**

NVLink data CRC Error Counter for Lane 10.

**#define**

**NVML\_FI\_DEV\_NVLINK\_CRC\_DATA\_ERROR\_COUNT\_L11  
107**

NVLink data CRC Error Counter for Lane 11.

**#define**

**NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_L6 108**

NVLink Replay Error Counter for Lane 6.

**#define**

**NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_L7 109**

NVLink Replay Error Counter for Lane 7.

**#define**

**NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_L8 110**

NVLink Replay Error Counter for Lane 8.

**#define**

**NVML\_FI\_DEV\_NVLINK\_REPLAY\_ERROR\_COUNT\_L9 111**

NVLink Replay Error Counter for Lane 9.

```
#define  
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L10 112  
NVLink Replay Error Counter for Lane 10.
```

```
#define  
NVML_FI_DEV_NVLINK_REPLAY_ERROR_COUNT_L11 113  
NVLink Replay Error Counter for Lane 11.
```

```
#define  
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L6  
114  
NVLink Recovery Error Counter for Lane 6.
```

```
#define  
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L7  
115  
NVLink Recovery Error Counter for Lane 7.
```

```
#define  
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L8  
116  
NVLink Recovery Error Counter for Lane 8.
```

```
#define  
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L9  
117  
NVLink Recovery Error Counter for Lane 9.
```

```
#define  
NVML_FI_DEV_NVLINK_RECOVERY_ERROR_COUNT_L10  
118  
NVLink Recovery Error Counter for Lane 10.
```

**#define**

**NVML\_FI\_DEV\_NVLINK\_RECOVERY\_ERROR\_COUNT\_L11**  
**119**

NVLink Recovery Error Counter for Lane 11.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C0\_L6 120**

NVLink Bandwidth Counter for Counter Set 0, Lane 6.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C0\_L7 121**

NVLink Bandwidth Counter for Counter Set 0, Lane 7.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C0\_L8 122**

NVLink Bandwidth Counter for Counter Set 0, Lane 8.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C0\_L9 123**

NVLink Bandwidth Counter for Counter Set 0, Lane 9.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C0\_L10 124**

NVLink Bandwidth Counter for Counter Set 0, Lane 10.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C0\_L11 125**

NVLink Bandwidth Counter for Counter Set 0, Lane 11.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C1\_L6 126**

NVLink Bandwidth Counter for Counter Set 1, Lane 6.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C1\_L7 127**

NVLink Bandwidth Counter for Counter Set 1, Lane 7.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C1\_L8 128**

NVLink Bandwidth Counter for Counter Set 1, Lane 8.

**#define NVML\_FI\_DEV\_NVLINK\_BANDWIDTH\_C1\_L9 129**

NVLink Bandwidth Counter for Counter Set 1, Lane 9.

```
#define NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L10 130
```

NVLink Bandwidth Counter for Counter Set 1, Lane 10.

```
#define NVML_FI_DEV_NVLINK_BANDWIDTH_C1_L11 131
```

NVLink Bandwidth Counter for Counter Set 1, Lane 11.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L6 132
```

NVLink Speed in MBps for Link 6.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L7 133
```

NVLink Speed in MBps for Link 7.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L8 134
```

NVLink Speed in MBps for Link 8.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L9 135
```

NVLink Speed in MBps for Link 9.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L10 136
```

NVLink Speed in MBps for Link 10.

```
#define NVML_FI_DEV_NVLINK_SPEED_MBPS_L11 137
```

NVLink Speed in MBps for Link 11.

```
#define NVML_FI_DEV_NVLINK_THROUGHPUT_DATA_TX  
138
```

NVLink TX Data throughput in KiB.

NVLink throughput counters field values

Link ID needs to be specified in the `scopeId` field in `nvmlFieldValue_t`. A `scopeId` of `UINT_MAX` returns aggregate value summed up across all links for the specified counter type in `fieldId`.

```
#define NVML_FI_DEV_NVLINK_THROUGHPUT_DATA_RX  
139
```

NVLink RX Data throughput in KiB.

```
#define NVML_FI_DEV_NVLINK_THROUGHPUT_RAW_TX  
140
```

NVLink TX Data + protocol overhead in KiB.

```
#define NVML_FI_DEV_NVLINK_THROUGHPUT_RAW_RX  
141
```

NVLink RX Data + protocol overhead in KiB.

```
#define NVML_FI_DEV_REMAPPED_COR 142
```

Number of remapped rows due to correctable errors.

```
#define NVML_FI_DEV_REMAPPED_UNC 143
```

Number of remapped rows due to uncorrectable errors.

```
#define NVML_FI_DEV_REMAPPED_PENDING 144
```

If any rows are pending remapping. 1=yes 0=no.

```
#define NVML_FI_DEV_REMAPPED_FAILURE 145
```

If any rows failed to be remapped 1=yes 0=no.

```
#define NVML_FI_DEV_NVLINK_REMOTE_NVLINK_ID 146
```

Remote device NVLink ID.

Remote device NVLink ID

Link ID needs to be specified in the scopeId field in `nvmlFieldValue_t`.

```
#define
```

```
NVML_FI_DEV_NVSWITCH_CONNECTED_LINK_COUNT 147
```

Number of NVLinks connected to NVSwitch.

NVSwitch: connected NVLink count

**#define**

**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L0 148**

NVLink data ECC Error Counter for Link 0.

**#define**

**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L1 149**

NVLink data ECC Error Counter for Link 1.

**#define**

**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L2 150**

NVLink data ECC Error Counter for Link 2.

**#define**

**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L3 151**

NVLink data ECC Error Counter for Link 3.

**#define**

**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L4 152**

NVLink data ECC Error Counter for Link 4.

**#define**

**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L5 153**

NVLink data ECC Error Counter for Link 5.

**#define**

**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L6 154**

NVLink data ECC Error Counter for Link 6.

**#define**

**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L7 155**

NVLink data ECC Error Counter for Link 7.



**#define**  
**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L8 156**  
NVLink data ECC Error Counter for Link 8.

**#define**  
**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L9 157**  
NVLink data ECC Error Counter for Link 9.

**#define**  
**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L10**  
**158**  
NVLink data ECC Error Counter for Link 10.

**#define**  
**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_L11**  
**159**  
NVLink data ECC Error Counter for Link 11.

**#define**  
**NVML\_FI\_DEV\_NVLINK\_ECC\_DATA\_ERROR\_COUNT\_TOTAL**  
**160**  
NVLink data ECC Error Counter total for all Links.

**#define NVML\_FI\_DEV\_NVLINK\_ERROR\_DL\_REPLAY 161**  
NVLink Replay Error Counter.

**#define NVML\_FI\_DEV\_NVLINK\_ERROR\_DL\_RECOVERY**  
**162**  
NVLink Recovery Error Counter.

**#define NVML\_FI\_DEV\_NVLINK\_ERROR\_DL\_CRC 163**  
NVLink CRC Error Counter.

**#define NVML\_FI\_DEV\_NVLINK\_GET\_SPEED 164**

NVLink Speed in MBps.

**#define NVML\_FI\_DEV\_NVLINK\_GET\_STATE 165**

NVLink State - Active,Inactive.

**#define NVML\_FI\_DEV\_NVLINK\_GET\_VERSION 166**

NVLink Version.

**#define NVML\_FI\_DEV\_NVLINK\_GET\_POWER\_STATE 167**

NVLink Power state. 0=HIGH\_SPEED 1=LOW\_SPEED.

**#define NVML\_FI\_DEV\_NVLINK\_GET\_POWER\_THRESHOLD  
168**

NVLink length of idle period (in units of 100us) before transitioning links to sleep state.

**#define NVML\_FI\_DEV\_PCIE\_L0\_TO\_RECOVERY\_COUNTER  
169**

Device PEX error recovery counter.

**#define NVML\_FI\_DEV\_C2C\_LINK\_COUNT 170**

Number of C2C Links present on the device.

**#define NVML\_FI\_DEV\_C2C\_LINK\_GET\_STATUS 171**

C2C Link Status 0=INACTIVE 1=ACTIVE.

**#define NVML\_FI\_DEV\_C2C\_LINK\_GET\_MAX\_BW 172**

C2C Link Speed in MBps for active links.

**#define NVML\_FI\_DEV\_POWER\_AVERAGE 185**

GPU power averaged over 1 sec interval, supported on Ampere (except GA100) or newer architectures.

Retrieves power usage for this GPU in milliwatts. It is only available if power management mode is supported. See [nvmlDeviceGetPowerManagementMode](#) and [nvmlDeviceGetPowerUsage](#).

scopeId needs to be specified. It signifies: 0 - GPU Only Scope - Metrics for GPU are retrieved 1 - Module scope - Metrics for the module (e.g. CPU + GPU) are retrieved. Note: CPU here refers to NVIDIA CPU (e.g. Grace). x86 or non-NVIDIA ARM is not supported

### **#define NVML\_FI\_DEV\_POWER\_INSTANT 186**

Current GPU power, supported on all architectures.

### **#define NVML\_FI\_DEV\_POWER\_MIN\_LIMIT 187**

Minimum power limit in milliwatts.

### **#define NVML\_FI\_DEV\_POWER\_MAX\_LIMIT 188**

Maximum power limit in milliwatts.

### **#define NVML\_FI\_DEV\_POWER\_DEFAULT\_LIMIT 189**

Default power limit in milliwatts (limit which device boots with).

### **#define NVML\_FI\_DEV\_POWER\_CURRENT\_LIMIT 190**

Limit currently enforced in milliwatts (This includes other limits set elsewhere. E.g. Out-of-band).

### **#define NVML\_FI\_DEV\_ENERGY 191**

Total energy consumption (in mJ) since the driver was last reloaded. Same as NVML\_FI\_DEV\_TOTAL\_ENERGY\_CONSUMPTION for the GPU.

### **#define NVML\_FI\_DEV\_POWER\_REQUESTED\_LIMIT 192**

Power limit requested by NVML or any other userspace client.

### **#define**

### **NVML\_FI\_DEV\_TEMPERATURE\_SHUTDOWN\_TLIMIT 193**

T.Limit temperature after which GPU may shut down for HW protection.

GPU T.Limit temperature thresholds in degree Celsius

These fields are supported on Ada and later architectures and supersedes `nvmlDeviceGetTemperatureThreshold`.

```
#define
```

```
NVML_FI_DEV_TEMPERATURE_SLOWDOWN_TLIMIT 194
```

T.Limit temperature after which GPU may begin HW slowdown.

```
#define NVML_FI_DEV_TEMPERATURE_MEM_MAX_TLIMIT  
195
```

T.Limit temperature after which GPU may begin SW slowdown due to memory temperature.

```
#define NVML_FI_DEV_TEMPERATURE_GPU_MAX_TLIMIT  
196
```

T.Limit temperature after which GPU may be throttled below base clock.

```
#define
```

```
NVML_FI_DEV_IS_MIG_MODE_INDEPENDENT_MIG_QUERY_CAPABLE  
199
```

MIG mode independent, MIG query capable device. 1=yes. 0=no.

```
#define NVML_FI_MAX 200
```

One greater than the largest field ID defined above.

## 4.4. Unit Structs

struct nvmlHwbcEntry\_t

struct nvmlLedState\_t

struct nvmlUnitInfo\_t

struct nvmlPSUInfo\_t

struct nvmlUnitFanInfo\_t

struct nvmlUnitFanSpeeds\_t

enum nvmlFanState\_t

Fan state enum.

#### Values

**NVML\_FAN\_NORMAL = 0**

Fan is working properly.

**NVML\_FAN\_FAILED = 1**

Fan has failed.

enum nvmlLedColor\_t

Led color enum.

#### Values

**NVML\_LED\_COLOR\_GREEN = 0**

GREEN, indicates good health.

**NVML\_LED\_COLOR\_AMBER = 1**

AMBER, indicates problem.

## 4.5. Accounting Statistics

Set of APIs designed to provide per process information about usage of GPU.



- ▶ All accounting statistics and accounting mode live in nvidia driver and reset to default (Disabled) when driver unloads. It is advised to run with persistence mode enabled.

- ▶ Enabling accounting mode has no negative impact on the GPU performance.

## struct nvmlAccountingStats\_t

### nvmlReturn\_t nvmlDeviceGetAccountingMode (nvmlDevice\_t device, nvmlEnableState\_t \*mode)

#### Parameters

##### device

The identifier of the target device

##### mode

Reference in which to return the current accounting mode

#### Returns

- ▶ NVML\_SUCCESS if the mode has been successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or mode are NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

#### Description

Queries the state of per process accounting mode.

For Kepler or newer fully supported devices.

See [nvmlDeviceGetAccountingStats](#) for more details. See [nvmlDeviceSetAccountingMode](#)

### nvmlReturn\_t nvmlDeviceGetAccountingStats (nvmlDevice\_t device, unsigned int pid, nvmlAccountingStats\_t \*stats)

#### Parameters

##### device

The identifier of the target device

##### pid

Process Id of the target process to query stats for

##### stats

Reference in which to return the process's accounting stats

## Returns

- ▶ NVML\_SUCCESS if stats have been successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or stats are NULL
- ▶ NVML\_ERROR\_NOT\_FOUND if process stats were not found
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if device doesn't support this feature or accounting mode is disabled or on vGPU host.
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

## Description

Queries process's accounting stats.

For Kepler or newer fully supported devices.

Accounting stats capture GPU utilization and other statistics across the lifetime of a process. Accounting stats can be queried during life time of the process and after its termination. The time field in `nvmlAccountingStats_t` is reported as 0 during the lifetime of the process and updated to actual running time after its termination. Accounting stats are kept in a circular buffer, newly created processes overwrite information about old processes.

See `nvmlAccountingStats_t` for description of each returned metric. List of processes that can be queried can be retrieved from `nvmlDeviceGetAccountingPids`.



- ▶ Accounting Mode needs to be on. See `nvmlDeviceGetAccountingMode`.
- ▶ Only compute and graphics applications stats can be queried. Monitoring applications stats can't be queried since they don't contribute to GPU utilization.
- ▶ In case of pid collision stats of only the latest process (that terminated last) will be reported

## See also:

`nvmlDeviceGetAccountingBufferSize`

`nvmlReturn_t nvmlDeviceGetAccountingPids`  
(`nvmlDevice_t device`, unsigned int \*count, unsigned int \*pids)

## Parameters

### device

The identifier of the target device

**count**

Reference in which to provide the pids array size, and to return the number of elements ready to be queried

**pids**

Reference in which to return list of process ids

**Returns**

- ▶ NVML\_SUCCESS if pids were successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or count is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if device doesn't support this feature or accounting mode is disabled or on vGPU host.
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if count is too small (count is set to expected value)
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Queries list of processes that can be queried for accounting stats. The list of processes returned can be in running or terminated state.

For Kepler or newer fully supported devices.

To just query the number of processes ready to be queried, call this function with \*count = 0 and pids=NULL. The return code will be NVML\_ERROR\_INSUFFICIENT\_SIZE, or NVML\_SUCCESS if list is empty.

For more details see [nvmlDeviceGetAccountingStats](#).



In case of PID collision some processes might not be accessible before the circular buffer is full.

**See also:**

[nvmlDeviceGetAccountingBufferSize](#)

**`nvmlReturn_t nvmlDeviceGetAccountingBufferSize  
(nvmlDevice_t device, unsigned int *bufferSize)`**

**Parameters****device**

The identifier of the target device



**bufferSize**

Reference in which to provide the size (in number of elements) of the circular buffer for accounting stats.

**Returns**

- ▶ NVML\_SUCCESS if buffer size was successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or bufferSize is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature or accounting mode is disabled
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Returns the number of processes that the circular buffer with accounting pids can hold.

For Kepler or newer fully supported devices.

This is the maximum number of processes that accounting information will be stored for before information about oldest processes will get overwritten by information about new processes.

**See also:**

[nvmlDeviceGetAccountingStats](#)

[nvmlDeviceGetAccountingPids](#)

## **`nvmlReturn_t nvmlDeviceSetAccountingMode` (`nvmlDevice_t device`, `nvmlEnableState_t mode`)**

**Parameters****device**

The identifier of the target device

**mode**

The target accounting mode

**Returns**

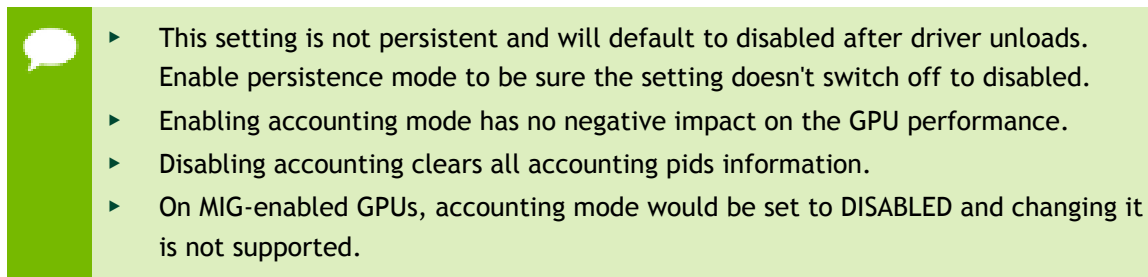
- ▶ NVML\_SUCCESS if the new mode has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device or mode are invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature

- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Enables or disables per process accounting.

For Kepler or newer fully supported devices. Requires root/admin permissions.



- ▶ This setting is not persistent and will default to disabled after driver unloads. Enable persistence mode to be sure the setting doesn't switch off to disabled.
- ▶ Enabling accounting mode has no negative impact on the GPU performance.
- ▶ Disabling accounting clears all accounting pids information.
- ▶ On MIG-enabled GPUs, accounting mode would be set to `DISABLED` and changing it is not supported.

See [nvmlDeviceGetAccountingMode](#) See [nvmlDeviceGetAccountingStats](#) See [nvmlDeviceClearAccountingPids](#)

## `nvmlReturn_t nvmlDeviceClearAccountingPids` (`nvmlDevice_t device`)

### Parameters

#### `device`

The identifier of the target device

### Returns

- ▶ `NVML_SUCCESS` if accounting information has been cleared
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device are invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device doesn't support this feature
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Clears accounting information about all processes that have already terminated.

For Kepler or newer fully supported devices. Requires root/admin permissions.

See [nvmlDeviceGetAccountingMode](#) See [nvmlDeviceGetAccountingStats](#) See [nvmlDeviceSetAccountingMode](#)

## 4.6. Encoder Structs

`struct nvmlEncoderSessionInfo_t`

`enum nvmlEncoderType_t`

Represents type of encoder for capacity can be queried

### Values

`NVML_ENCODER_QUERY_H264 = 0x00`

H264 encoder.

`NVML_ENCODER_QUERY_HEVC = 0x01`

HEVC encoder.

`NVML_ENCODER_QUERY_AV1 = 0x02`

AV1 encoder.

`NVML_ENCODER_QUERY_UNKNOWN = 0xFF`

Unknown encoder.

## 4.7. Frame Buffer Capture Structures

`struct nvmlFBCStats_t`

`struct nvmlFBCSessionInfo_t`

`enum nvmlFBCSessionType_t`

Represents frame buffer capture session type

### Values

`NVML_FBC_SESSION_TYPE_UNKNOWN = 0`

Unknown.

`NVML_FBC_SESSION_TYPE_TOSYS`

ToSys.

NVML\_FBC\_SESSION\_TYPE\_CUDA

Cuda.

NVML\_FBC\_SESSION\_TYPE\_VID

Vid.

NVML\_FBC\_SESSION\_TYPE\_HWENC

HEnc.

```
#define NVML_NVFBC_SESSION_FLAG_DIFFMAP_ENABLED  
0x00000001
```

Bit specifying differential map state.

```
#define  
NVML_NVFBC_SESSION_FLAG_CLASSIFICATIONMAP_ENABLED  
0x00000002
```

Bit specifying classification map state.

```
#define  
NVML_NVFBC_SESSION_FLAG_CAPTURE_WITH_WAIT_NO_WAIT  
0x00000004
```

Bit specifying if capture was requested as non-blocking call.

```
#define  
NVML_NVFBC_SESSION_FLAG_CAPTURE_WITH_WAIT_INFINITE  
0x00000008
```

Bit specifying if capture was requested as blocking call.

```
#define  
NVML_NVFBC_SESSION_FLAG_CAPTURE_WITH_WAIT_TIMEOUT  
0x00000010
```

Bit specifying if capture was requested as blocking call with timeout period.

## 4.8. definitions related to the drain state

## enum nvmlDetachGpuState\_t

Is the GPU device to be removed from the kernel by nvmlDeviceRemoveGpu()

### Values

NVML\_DETACH\_GPU\_KEEP = 0  
 NVML\_DETACH\_GPU\_REMOVE

## enum nvmlPcieLinkState\_t

Parent bridge PCIe link state requested by nvmlDeviceRemoveGpu()

### Values

NVML\_PCIE\_LINK\_KEEP = 0  
 NVML\_PCIE\_LINK\_SHUT\_DOWN

## 4.9. /nvmlDevice definitions related to Confidential Computing

struct nvmlSystemConfComputeSettings\_v1\_t

struct nvmlConfComputeMemSizeInfo\_t

#define NVML\_CC\_SYSTEM\_CPU\_CAPS\_NONE 0

Confidential Compute CPU Capabilities values

#define NVML\_CC\_SYSTEM\_GPUS\_CC\_NOT\_CAPABLE 0

Confidential Compute GPU Capabilities values

#define NVML\_CC\_SYSTEM\_DEVTOOLS\_MODE\_OFF 0

Confidential Compute DevTools Mode values

#define NVML\_CC\_SYSTEM\_ENVIRONMENT\_UNAVAILABLE  
 0

Confidential Compute Environment values

```
#define NVML_CC_SYSTEM_FEATURE_DISABLED 0
```

Confidential Compute Feature Status values

```
#define NVML_CC_SYSTEM_MULTIGPU_NONE 0
```

Confidential Compute Multigpu mode values

```
#define NVML_CC_ACCEPTING_CLIENT_REQUESTS_FALSE
0
```

Confidential Compute GPUs/System Ready State values

```
#define NVML_GPU_CERT_CHAIN_SIZE 0x1000
```

GPU Certificate Details

```
#define NVML_CC_GPU_CEC_NONCE_SIZE 0x20
```

GPU Attestation Report

## 4.10. Initialization and Cleanup

This chapter describes the methods that handle NVML initialization and cleanup. It is the user's responsibility to call `nvmlInit_v2()` before calling any other methods, and `nvmlShutdown()` once NVML is no longer being used.

### `nvmlReturn_t nvmlInit_v2 (void)`

#### Returns

- ▶ `NVML_SUCCESS` if NVML has been properly initialized
- ▶ `NVML_ERROR_DRIVER_NOT_LOADED` if NVIDIA driver is not running
- ▶ `NVML_ERROR_NO_PERMISSION` if NVML does not have permission to talk to the driver
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

#### Description

Initialize NVML, but don't initialize any GPUs yet.



- ▶ `nvmlInit_v3` introduces a "flags" argument, that allows passing boolean values modifying the behaviour of `nvmlInit()`.
- ▶ In NVML 5.319 new `nvmlInit_v2` has replaced `nvmlInit_v1` (default in NVML 4.304 and older) that did initialize all GPU devices in the system.

This allows NVML to communicate with a GPU when other GPUs in the system are unstable or in a bad state. When using this API, GPUs are discovered and initialized in `nvmlDeviceGetHandleBy*` functions instead.



To contrast `nvmlInit_v2` with `nvmlInit_v1`, NVML 4.304 `nvmlInit_v1` will fail when any detected GPU is in a bad or unstable state.

For all products.

This method, should be called once before invoking any other methods in the library. A reference count of the number of initializations is maintained. Shutdown only occurs when the reference count reaches zero.

## `nvmlReturn_t nvmlInitWithFlags (unsigned int flags)`

### Parameters

#### **flags**

behaviour modifier flags

### Returns

- ▶ `NVML_SUCCESS` if NVML has been properly initialized
- ▶ `NVML_ERROR_DRIVER_NOT_LOADED` if NVIDIA driver is not running
- ▶ `NVML_ERROR_NO_PERMISSION` if NVML does not have permission to talk to the driver
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

`nvmlInitWithFlags` is a variant of `nvmlInit()`, that allows passing a set of boolean values modifying the behaviour of `nvmlInit()`. Other than the "flags" parameter it is completely similar to [`nvmlInit\_v2`](#).

For all products.

## `nvmlReturn_t nvmlShutdown (void)`

### Returns

- ▶ `NVML_SUCCESS` if NVML has been properly shut down
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Shut down NVML by releasing all GPU resources previously allocated with `nvmlInit_v2()`.

For all products.

This method should be called after NVML work is done, once for each call to `nvmlInit_v2()`. A reference count of the number of initializations is maintained. Shutdown only occurs when the reference count reaches zero. For backwards compatibility, no error is reported if `nvmlShutdown()` is called more times than `nvmlInit()`.

### `#define NVML_INIT_FLAG_NO_GPUS 1`

Don't fail `nvmlInit()` when no GPUs are found.

### `#define NVML_INIT_FLAG_NO_ATTACH 2`

Don't attach GPUs.

## 4.11. Error reporting

This chapter describes helper functions for error reporting routines.

### `const DECLDIR char *nvmlErrorString (nvmlReturn_t result)`

#### Parameters

##### **result**

NVML error code to convert

#### Returns

String representation of the error.



## Description

Helper method for converting NVML error codes into readable strings.

For all products.

## 4.12. Constants

```
#define NVML_DEVICE_INFOROM_VERSION_BUFFER_SIZE  
16
```

Buffer size guaranteed to be large enough for `nvmDeviceGetInforomVersion` and `nvmDeviceGetInforomImageVersion`

```
#define NVML_DEVICE_UUID_BUFFER_SIZE 80
```

Buffer size guaranteed to be large enough for storing GPU identifiers.

```
#define NVML_DEVICE_UUID_V2_BUFFER_SIZE 96
```

Buffer size guaranteed to be large enough for `nvmDeviceGetUUID`

```
#define NVML_DEVICE_PART_NUMBER_BUFFER_SIZE 80
```

Buffer size guaranteed to be large enough for `nvmDeviceGetBoardPartNumber`

```
#define NVML_SYSTEM_DRIVER_VERSION_BUFFER_SIZE  
80
```

Buffer size guaranteed to be large enough for `nvmSystemGetDriverVersion`

```
#define NVML_SYSTEM_NVML_VERSION_BUFFER_SIZE 80
```

Buffer size guaranteed to be large enough for `nvmSystemGetNVMLVersion`

```
#define NVML_DEVICE_NAME_BUFFER_SIZE 64
```

Buffer size guaranteed to be large enough for storing GPU device names.

```
#define NVML_DEVICE_NAME_V2_BUFFER_SIZE 96
```

Buffer size guaranteed to be large enough for `nvmDeviceGetName`

```
#define NVML_DEVICE_SERIAL_BUFFER_SIZE 30
```

Buffer size guaranteed to be large enough for `nvmIDeviceGetSerial`

```
#define NVML_DEVICE_VBIOS_VERSION_BUFFER_SIZE 32
```

Buffer size guaranteed to be large enough for `nvmIDeviceGetVbiosVersion`

## 4.13. System Queries

This chapter describes the queries that NVML can perform against the local system. These queries are not device-specific.

```
nvmIReturn_t nvmISystemGetDriverVersion (char  
*version, unsigned int length)
```

### Parameters

#### **version**

Reference in which to return the version identifier

#### **length**

The maximum allowed length of the string returned in version

### Returns

- ▶ NVML\_SUCCESS if version has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if version is NULL
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if length is too small

### Description

Retrieves the version of the system's graphics driver.

For all products.

The version identifier is an alphanumeric string. It will not exceed 80 characters in length (including the NULL terminator). See `nvmIConstants::NVML_SYSTEM_DRIVER_VERSION_BUFFER_SIZE`.

## `nvmlReturn_t nvmlSystemGetNVMLVersion (char *version, unsigned int length)`

### Parameters

#### **version**

Reference in which to return the version identifier

#### **length**

The maximum allowed length of the string returned in version

### Returns

- ▶ `NVML_SUCCESS` if version has been set
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if version is `NULL`
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if length is too small

### Description

Retrieves the version of the NVML library.

For all products.

The version identifier is an alphanumeric string. It will not exceed 80 characters in length (including the `NULL` terminator). See [nvmlConstants::NVML\\_SYSTEM\\_NVML\\_VERSION\\_BUFFER\\_SIZE](#).

## `nvmlReturn_t nvmlSystemGetCudaDriverVersion (int *cudaDriverVersion)`

### Parameters

#### **cudaDriverVersion**

Reference in which to return the version identifier

### Returns

- ▶ `NVML_SUCCESS` if `cudaDriverVersion` has been set
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `cudaDriverVersion` is `NULL`

### Description

Retrieves the version of the CUDA driver.

For all products.

The CUDA driver version returned will be retrieved from the currently installed version of CUDA. If the cuda library is not found, this function will return a known supported version number.

## **nvmlReturn\_t nvmlSystemGetCudaDriverVersion\_v2 (int \*cudaDriverVersion)**

### **Parameters**

#### **cudaDriverVersion**

Reference in which to return the version identifier

### **Returns**

- ▶ NVML\_SUCCESS if cudaDriverVersion has been set
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if cudaDriverVersion is NULL
- ▶ NVML\_ERROR\_LIBRARY\_NOT\_FOUND if libcuda.so.1 or libcuda.dll is not found
- ▶ NVML\_ERROR\_FUNCTION\_NOT\_FOUND if cuDriverGetVersion() is not found in the shared library

### **Description**

Retrieves the version of the CUDA driver from the shared library.

For all products.

The returned CUDA driver version by calling cuDriverGetVersion()

## **nvmlReturn\_t nvmlSystemGetProcessName (unsigned int pid, char \*name, unsigned int length)**

### **Parameters**

#### **pid**

The identifier of the process

#### **name**

Reference in which to return the process name

#### **length**

The maximum allowed length of the string returned in name

### **Returns**

- ▶ NVML\_SUCCESS if name has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized

- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if name is NULL or length is 0.
- ▶ NVML\_ERROR\_NOT\_FOUND if process doesn't exist
- ▶ NVML\_ERROR\_NO\_PERMISSION if the user doesn't have permission to perform this operation
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Gets name of the process with provided process id

For all products.

Returned process name is cropped to provided length. name string is encoded in ANSI.

**nvmlReturn\_t nvmlSystemGetHicVersion (unsigned int \*hwbcCount, nvmlHwbcEntry\_t \*hwbcEntries)**

### Parameters

#### hwbcCount

Size of hwbcEntries array

#### hwbcEntries

Array holding information about hwbc

### Returns

- ▶ NVML\_SUCCESS if hwbcCount and hwbcEntries have been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if either hwbcCount or hwbcEntries is NULL
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if hwbcCount indicates that the hwbcEntries array is too small

### Description

Retrieves the IDs and firmware versions for any Host Interface Cards (HICs) in the system.

For S-class products.

The hwbcCount argument is expected to be set to the size of the input hwbcEntries array. The HIC must be connected to an S-class system for it to be reported by this function.

**nvmlReturn\_t nvmlSystemGetTopologyGpuSet (unsigned int cpuNumber, unsigned int \*count, nvmlDevice\_t \*deviceArray)**

#### Parameters

##### **cpuNumber**

The CPU number

##### **count**

When zero, is set to the number of matching GPUs such that deviceArray can be malloc'd. When non-zero, deviceArray will be filled with count number of device handles.

##### **deviceArray**

An array of device handles for GPUs found with affinity to cpuNumber

#### Returns

- ▶ NVML\_SUCCESS if deviceArray or count (if initially zero) has been set
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if cpuNumber, or count is invalid, or deviceArray is NULL with a non-zero count
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device or OS does not support this feature
- ▶ NVML\_ERROR\_UNKNOWN an error has occurred in underlying topology discovery

#### Description

Retrieve the set of GPUs that have a CPU affinity with the given CPU number For all products. Supported on Linux only.

**#define NVML\_CUDA\_DRIVER\_VERSION\_MAJOR  
((v)/1000)**

Macros for converting the CUDA driver version number to Major and Minor version numbers.

## 4.14. Unit Queries

This chapter describes that queries that NVML can perform against each unit. For S-class systems only. In each case the device is identified with an nvmlUnit\_t handle. This handle is obtained by calling `nvmlUnitGetHandleByIndex()`.

## `nvmlReturn_t nvmlUnitGetCount (unsigned int *unitCount)`

### Parameters

#### **unitCount**

Reference in which to return the number of units

### Returns

- ▶ `NVML_SUCCESS` if `unitCount` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `unitCount` is `NULL`
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the number of units in the system.

For S-class products.

## `nvmlReturn_t nvmlUnitGetHandleByIndex (unsigned int index, nvmlUnit_t *unit)`

### Parameters

#### **index**

The index of the target unit,  $\geq 0$  and  $< \text{unitCount}$

#### **unit**

Reference in which to return the unit handle

### Returns

- ▶ `NVML_SUCCESS` if `unit` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `index` is invalid or `unit` is `NULL`
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Acquire the handle for a particular unit, based on its index.

For S-class products.

Valid indices are derived from the `unitCount` returned by `nvmlUnitGetCount()`. For example, if `unitCount` is 2 the valid indices are 0 and 1, corresponding to UNIT 0 and UNIT 1.

The order in which NVML enumerates units has no guarantees of consistency between reboots.

## `nvmlReturn_t nvmlUnitGetUnitInfo (nvmlUnit_t unit, nvmlUnitInfo_t *info)`

### Parameters

#### `unit`

The identifier of the target unit

#### `info`

Reference in which to return the unit information

### Returns

- ▶ `NVML_SUCCESS` if `info` has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `unit` is invalid or `info` is `NULL`

### Description

Retrieves the static information associated with a unit.

For S-class products.

See `nvmlUnitInfo_t` for details on available unit info.

## `nvmlReturn_t nvmlUnitGetLedState (nvmlUnit_t unit, nvmlLedState_t *state)`

### Parameters

#### `unit`

The identifier of the target unit

#### `state`

Reference in which to return the current LED state

### Returns

- ▶ `NVML_SUCCESS` if `state` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized



- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if unit is invalid or state is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this is not an S-class product
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the LED state associated with this unit.

For S-class products.

See [nvmlLedState\\_t](#) for details on allowed states.

### See also:

[nvmlUnitSetLedState\(\)](#)

## **`nvmlReturn_t nvmlUnitGetPsuInfo (nvmlUnit_t unit, nvmlPSUInfo_t *psu)`**

### Parameters

#### **unit**

The identifier of the target unit

#### **psu**

Reference in which to return the PSU information

### Returns

- ▶ NVML\_SUCCESS if psu has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if unit is invalid or psu is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this is not an S-class product
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the PSU stats for the unit.

For S-class products.

See [nvmlPSUInfo\\_t](#) for details on available PSU info.

## `nvmlReturn_t nvmlUnitGetTemperature (nvmlUnit_t unit, unsigned int type, unsigned int *temp)`

### Parameters

#### **unit**

The identifier of the target unit

#### **type**

The type of reading to take

#### **temp**

Reference in which to return the intake temperature

### Returns

- ▶ `NVML_SUCCESS` if temp has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if unit or type is invalid or temp is NULL
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this is not an S-class product
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the temperature readings for the unit, in degrees C.

For S-class products.

Depending on the product, readings may be available for intake (type=0), exhaust (type=1) and board (type=2).

## `nvmlReturn_t nvmlUnitGetFanSpeedInfo (nvmlUnit_t unit, nvmlUnitFanSpeeds_t *fanSpeeds)`

### Parameters

#### **unit**

The identifier of the target unit

#### **fanSpeeds**

Reference in which to return the fan speed information

### Returns

- ▶ `NVML_SUCCESS` if fanSpeeds has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if unit is invalid or fanSpeeds is NULL

- ▶ `NVML_ERROR_NOT_SUPPORTED` if this is not an S-class product
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the fan speed readings for the unit.

For S-class products.

See `nvmlUnitFanSpeeds_t` for details on available fan speed info.

`nvmlReturn_t nvmlUnitGetDevices (nvmlUnit_t unit, unsigned int *deviceCount, nvmlDevice_t *devices)`

### Parameters

#### **unit**

The identifier of the target unit

#### **deviceCount**

Reference in which to provide the devices array size, and to return the number of attached GPU devices

#### **devices**

Reference in which to return the references to the attached GPU devices

### Returns

- ▶ `NVML_SUCCESS` if deviceCount and devices have been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if deviceCount indicates that the devices array is too small
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if unit is invalid, either of deviceCount or devices is NULL
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the set of GPU devices that are attached to the specified unit.

For S-class products.

The deviceCount argument is expected to be set to the size of the input devices array.

## 4.15. Device Queries

This chapter describes that queries that NVML can perform against each device. In each case the device is identified with an `nvmlDevice_t` handle. This handle is obtained by calling one of `nvmlDeviceGetHandleByIndex_v2()`, `nvmlDeviceGetHandleBySerial()`, `nvmlDeviceGetHandleByPciBusId_v2()`, or `nvmlDeviceGetHandleByUUID()`.

### CPU and Memory Affinity

#### `nvmlReturn_t nvmlDeviceGetCount_v2 (unsigned int *deviceCount)`

##### Parameters

##### `deviceCount`

Reference in which to return the number of accessible devices

##### Returns

- ▶ `NVML_SUCCESS` if `deviceCount` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `deviceCount` is `NULL`
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

##### Description

Retrieves the number of compute devices in the system. A compute device is a single GPU.

For all products.

Note: New `nvmlDeviceGetCount_v2` (default in NVML 5.319) returns count of all devices in the system even if `nvmlDeviceGetHandleByIndex_v2` returns `NVML_ERROR_NO_PERMISSION` for such device. Update your code to handle this error, or use NVML 4.304 or older `nvml` header file. For backward binary compatibility reasons `_v1` version of the API is still present in the shared library. Old `_v1` version of `nvmlDeviceGetCount` doesn't count devices that NVML has no permission to talk to.

## `nvmlReturn_t nvmlDeviceGetAttributes_v2` (`nvmlDevice_t device, nvmlDeviceAttributes_t *attributes`)

### Parameters

#### **device**

NVML device handle

#### **attributes**

Device attributes

### Returns

- ▶ `NVML_SUCCESS` if device attributes were successfully retrieved
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device handle is invalid
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Get attributes (engine counts etc.) for the given NVML device handle.



This API currently only supports MIG device handles.

For Ampere or newer fully supported devices. Supported on Linux only.

## `nvmlReturn_t nvmlDeviceGetHandleByIndex_v2` (`unsigned int index, nvmlDevice_t *device`)

### Parameters

#### **index**

The index of the target GPU,  $\geq 0$  and  $<$  `accessibleDevices`

#### **device**

Reference in which to return the device handle

### Returns

- ▶ `NVML_SUCCESS` if device has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if index is invalid or device is `NULL`

- ▶ `NVML_ERROR_INSUFFICIENT_POWER` if any attached devices have improperly attached external power cables
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to talk to this device
- ▶ `NVML_ERROR_IRQ_ISSUE` if NVIDIA kernel detected an interrupt issue with the attached GPUs
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Acquire the handle for a particular device, based on its index.

For all products.

Valid indices are derived from the `accessibleDevices` count returned by `nvmlDeviceGetCount_v2()`. For example, if `accessibleDevices` is 2 the valid indices are 0 and 1, corresponding to GPU 0 and GPU 1.

The order in which NVML enumerates devices has no guarantees of consistency between reboots. For that reason it is recommended that devices be looked up by their PCI ids or UUID. See `nvmlDeviceGetHandleByUUID()` and `nvmlDeviceGetHandleByPciBusId_v2()`.

Note: The NVML index may not correlate with other APIs, such as the CUDA device index.

Starting from NVML 5, this API causes NVML to initialize the target GPU NVML may initialize additional GPUs if:

- ▶ The target GPU is an SLI slave

Note: New `nvmlDeviceGetCount_v2` (default in NVML 5.319) returns count of all devices in the system even if `nvmlDeviceGetHandleByIndex_v2` returns `NVML_ERROR_NO_PERMISSION` for such device. Update your code to handle this error, or use NVML 4.304 or older `nvml` header file. For backward binary compatibility reasons `_v1` version of the API is still present in the shared library. Old `_v1` version of `nvmlDeviceGetCount` doesn't count devices that NVML has no permission to talk to.

This means that `nvmlDeviceGetHandleByIndex_v2` and `_v1` can return different devices for the same index. If you don't touch macros that map old (`_v1`) versions to `_v2` versions at the top of the file you don't need to worry about that.

### See also:

[nvmlDeviceGetIndex](#)

[nvmlDeviceGetCount](#)

## `nvmlReturn_t nvmlDeviceGetHandleBySerial (const char *serial, nvmlDevice_t *device)`

### Parameters

#### **serial**

The board serial number of the target GPU

#### **device**

Reference in which to return the device handle

### Returns

- ▶ `NVML_SUCCESS` if device has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if serial is invalid, device is `NULL` or more than one device has the same serial (dual GPU boards)
- ▶ `NVML_ERROR_NOT_FOUND` if serial does not match a valid device on the system
- ▶ `NVML_ERROR_INSUFFICIENT_POWER` if any attached devices have improperly attached external power cables
- ▶ `NVML_ERROR_IRQ_ISSUE` if NVIDIA kernel detected an interrupt issue with the attached GPUs
- ▶ `NVML_ERROR_GPU_IS_LOST` if any GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Acquire the handle for a particular device, based on its board serial number.

For Fermi or newer fully supported devices.

This number corresponds to the value printed directly on the board, and to the value returned by `nvmlDeviceGetSerial()`.

**Deprecated** Since more than one GPU can exist on a single board this function is deprecated in favor of `nvmlDeviceGetHandleByUUID`. For dual GPU boards this function will return `NVML_ERROR_INVALID_ARGUMENT`.

Starting from NVML 5, this API causes NVML to initialize the target GPU NVML may initialize additional GPUs as it searches for the target GPU

### See also:

[nvmlDeviceGetSerial](#)

[nvmlDeviceGetHandleByUUID](#)

## `nvmlReturn_t nvmlDeviceGetHandleByUUID (const char *uuid, nvmlDevice_t *device)`

### Parameters

#### **uuid**

The UUID of the target GPU or MIG instance

#### **device**

Reference in which to return the device handle or MIG device handle

### Returns

- ▶ `NVML_SUCCESS` if device has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `uuid` is invalid or `device` is null
- ▶ `NVML_ERROR_NOT_FOUND` if `uuid` does not match a valid device on the system
- ▶ `NVML_ERROR_INSUFFICIENT_POWER` if any attached devices have improperly attached external power cables
- ▶ `NVML_ERROR_IRQ_ISSUE` if NVIDIA kernel detected an interrupt issue with the attached GPUs
- ▶ `NVML_ERROR_GPU_IS_LOST` if any GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Acquire the handle for a particular device, based on its globally unique immutable UUID associated with each device.

For all products.

Starting from NVML 5, this API causes NVML to initialize the target GPU NVML may initialize additional GPUs as it searches for the target GPU

### See also:

[nvmlDeviceGetUUID](#)



## `nvmlReturn_t nvmlDeviceGetHandleByPciBusId_v2` (`const char *pciBusId, nvmlDevice_t *device`)

### Parameters

#### `pciBusId`

The PCI bus id of the target GPU

#### `device`

Reference in which to return the device handle

### Returns

- ▶ `NVML_SUCCESS` if device has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `pciBusId` is invalid or `device` is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if `pciBusId` does not match a valid device on the system
- ▶ `NVML_ERROR_INSUFFICIENT_POWER` if the attached device has improperly attached external power cables
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to talk to this device
- ▶ `NVML_ERROR_IRQ_ISSUE` if NVIDIA kernel detected an interrupt issue with the attached GPUs
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Acquire the handle for a particular device, based on its PCI bus id.

For all products.

This value corresponds to the `nvmlPciInfo_t::busId` returned by `nvmlDeviceGetPciInfo_v3()`.

Starting from NVML 5, this API causes NVML to initialize the target GPU NVML may initialize additional GPUs if:

- ▶ The target GPU is an SLI slave



NVML 4.304 and older version of `nvmlDeviceGetHandleByPciBusId_v1` returns `NVML_ERROR_NOT_FOUND` instead of `NVML_ERROR_NO_PERMISSION`.

## `nvmlReturn_t nvmlDeviceGetName (nvmlDevice_t device, char *name, unsigned int length)`

### Parameters

#### **device**

The identifier of the target device

#### **name**

Reference in which to return the product name

#### **length**

The maximum allowed length of the string returned in name

### Returns

- ▶ `NVML_SUCCESS` if name has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, or name is `NULL`
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if length is too small
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the name of this device.

For all products.

The name is an alphanumeric string that denotes a particular product, e.g. Tesla C2070. It will not exceed 96 characters in length (including the `NULL` terminator). See [nvmlConstants::NVML\\_DEVICE\\_NAME\\_V2\\_BUFFER\\_SIZE](#).

When used with MIG device handles the API returns MIG device names which can be used to identify devices based on their attributes.

## `nvmlReturn_t nvmlDeviceGetBrand (nvmlDevice_t device, nvmlBrandType_t *type)`

### Parameters

#### **device**

The identifier of the target device

#### **type**

Reference in which to return the product brand type

**Returns**

- ▶ NVML\_SUCCESS if name has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, or type is NULL
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the brand of this device.

For all products.

The type is a member of `nvmlBrandType_t` defined above.

## `nvmlReturn_t nvmlDeviceGetIndex (nvmlDevice_t device, unsigned int *index)`

**Parameters****device**

The identifier of the target device

**index**

Reference in which to return the NVML index of the device

**Returns**

- ▶ NVML\_SUCCESS if index has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, or index is NULL
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the NVML index of this device.

For all products.

Valid indices are derived from the `accessibleDevices` count returned by `nvmlDeviceGetCount_v2()`. For example, if `accessibleDevices` is 2 the valid indices are 0 and 1, corresponding to GPU 0 and GPU 1.

The order in which NVML enumerates devices has no guarantees of consistency between reboots. For that reason it is recommended that devices be looked up by their PCI ids or GPU UUID. See [nvmlDeviceGetHandleByPciBusId\\_v2\(\)](#) and [nvmlDeviceGetHandleByUUID\(\)](#).

When used with MIG device handles this API returns indices that can be passed to [nvmlDeviceGetMigDeviceHandleByIndex](#) to retrieve an identical handle. MIG device indices are unique within a device.

Note: The NVML index may not correlate with other APIs, such as the CUDA device index.

**See also:**

[nvmlDeviceGetHandleByIndex\(\)](#)

[nvmlDeviceGetCount\(\)](#)

## [nvmlReturn\\_t nvmlDeviceGetSerial \(nvmlDevice\\_t device, char \\*serial, unsigned int length\)](#)

### Parameters

**device**

The identifier of the target device

**serial**

Reference in which to return the board/module serial number

**length**

The maximum allowed length of the string returned in serial

### Returns

- ▶ NVML\_SUCCESS if serial has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, or serial is NULL
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if length is too small
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the globally unique board serial number associated with this device's board.

For all products with an inforom.

The serial number is an alphanumeric string that will not exceed 30 characters (including the NULL terminator). This number matches the serial number tag that is physically attached to the board. See `nvmlConstants::NVML_DEVICE_SERIAL_BUFFER_SIZE`.

## `nvmlReturn_t nvmlDeviceGetC2cModeInfoV (nvmlDevice_t device, nvmlC2cModeInfo_v1_t *c2cModeInfo)`

### Parameters

#### `device`

The identifier of the target device

#### `c2cModeInfo`

Output struct containing the device's C2C Mode info

### Returns

- ▶ `NVML_SUCCESS` if C2C Mode Infor query is successful
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, or serial is NULL
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the Device's C2C Mode information

## `nvmlReturn_t nvmlDeviceGetTopologyCommonAncestor (nvmlDevice_t device1, nvmlDevice_t device2, nvmlGpuTopologyLevel_t *pathInfo)`

### Parameters

#### `device1`

The identifier of the first device

#### `device2`

The identifier of the second device

#### `pathInfo`

A `nvmlGpuTopologyLevel_t` that gives the path type

**Returns**

- ▶ NVML\_SUCCESS if pathInfo has been set
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device1, or device2 is invalid, or pathInfo is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device or OS does not support this feature
- ▶ NVML\_ERROR\_UNKNOWN an error has occurred in underlying topology discovery

**Description**

Retrieve the common ancestor for two devices For all products. Supported on Linux only.

**nvmlReturn\_t nvmlDeviceGetTopologyNearestGpus**  
**(nvmlDevice\_t device, nvmlGpuTopologyLevel\_t level,**  
**unsigned int \*count, nvmlDevice\_t \*deviceArray)**

**Parameters****device**

The identifier of the first device

**level**

The `nvmlGpuTopologyLevel_t` level to search for other GPUs

**count**

When zero, is set to the number of matching GPUs such that deviceArray can be malloc'd. When non-zero, deviceArray will be filled with count number of device handles.

**deviceArray**

An array of device handles for GPUs found at level

**Returns**

- ▶ NVML\_SUCCESS if deviceArray or count (if initially zero) has been set
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, level, or count is invalid, or deviceArray is NULL with a non-zero count
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device or OS does not support this feature
- ▶ NVML\_ERROR\_UNKNOWN an error has occurred in underlying topology discovery

**Description**

Retrieve the set of GPUs that are nearest to a given device at a specific interconnectivity level For all products. Supported on Linux only.

```
nvmlReturn_t nvmlDeviceGetP2PStatus
(nvmlDevice_t device1, nvmlDevice_t device2,
nvmlGpuP2PCapsIndex_t p2pIndex, nvmlGpuP2PStatus_t
*p2pStatus)
```

**Parameters****device1**

The first device

**device2**

The second device

**p2pIndex**

p2p Capability Index being looked for between device1 and device2

**p2pStatus**

Reference in which to return the status of the p2pIndex between device1 and device2

**Returns**

- ▶ NVML\_SUCCESS if p2pStatus has been populated
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device1 or device2 or p2pIndex is invalid or p2pStatus is NULL
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve the status for a given p2p capability index between a given pair of GPU

```
nvmlReturn_t nvmlDeviceGetUUID (nvmlDevice_t
device, char *uuid, unsigned int length)
```

**Parameters****device**

The identifier of the target device

**uuid**

Reference in which to return the GPU UUID

**length**

The maximum allowed length of the string returned in `uuid`

**Returns**

- ▶ `NVML_SUCCESS` if `uuid` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, or `uuid` is `NULL`
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `length` is too small
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

**Description**

Retrieves the globally unique immutable UUID associated with this device, as a 5 part hexadecimal string, that augments the immutable, board serial identifier.

For all products.

The UUID is a globally unique identifier. It is the only available identifier for pre-Fermi-architecture products. It does NOT correspond to any identifier printed on the board. It will not exceed 96 characters in length (including the `NULL` terminator). See [`nvmlConstants::NVML\_DEVICE\_UUID\_V2\_BUFFER\_SIZE`](#).

When used with MIG device handles the API returns globally unique UUIDs which can be used to identify MIG devices across both GPU and MIG devices. UUIDs are immutable for the lifetime of a MIG device.

## `nvmlReturn_t nvmlDeviceGetMinorNumber` (`nvmlDevice_t device`, `unsigned int *minorNumber`)

**Parameters****device**

The identifier of the target device

**minorNumber**

Reference in which to return the minor number for the device

**Returns**

- ▶ `NVML_SUCCESS` if the minor number is successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized



- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or minorNumber is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves minor number for the device. The minor number for the device is such that the Nvidia device node file for each GPU will have the form /dev/nvidia[minor number].

For all products. Supported only for Linux

## nvmlReturn\_t nvmlDeviceGetBoardPartNumber (nvmlDevice\_t device, char \*partNumber, unsigned int length)

### Parameters

#### device

Identifier of the target device

#### partNumber

Reference to the buffer to return

#### length

Length of the buffer reference

### Returns

- ▶ NVML\_SUCCESS if partNumber has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the needed VBIOS fields have not been filled
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or serial is NULL
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the the device board part number which is programmed into the board's InfoROM

For all products.

## `nvmlReturn_t nvmlDeviceGetInforomVersion` (`nvmlDevice_t device`, `nvmlInforomObject_t object`, `char *version`, `unsigned int length`)

### Parameters

#### **device**

The identifier of the target device

#### **object**

The target infoROM object

#### **version**

Reference in which to return the infoROM version

#### **length**

The maximum allowed length of the string returned in version

### Returns

- ▶ `NVML_SUCCESS` if version has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if version is `NULL`
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if length is too small
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not have an infoROM
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the version information for the device's infoROM object.

For all products with an inforom.

Fermi and higher parts have non-volatile on-board memory for persisting device info, such as aggregate ECC counts. The version of the data structures in this memory may change from time to time. It will not exceed 16 characters in length (including the `NULL` terminator). See `nvmlConstants::NVML_DEVICE_INFOROM_VERSION_BUFFER_SIZE`.

See `nvmlInforomObject_t` for details on the available infoROM objects.

### See also:

[nvmlDeviceGetInforomImageVersion](#)

## `nvmReturn_t nvmDeviceGetInforomImageVersion` (`nvmDevice_t device`, `char *version`, `unsigned int length`)

### Parameters

**device**

The identifier of the target device

**version**

Reference in which to return the infoROM image version

**length**

The maximum allowed length of the string returned in version

### Returns

- ▶ `NVML_SUCCESS` if version has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if version is `NULL`
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if length is too small
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not have an infoROM
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the global infoROM image version

For all products with an inforom.

Image version just like VBIOS version uniquely describes the exact version of the infoROM flashed on the board in contrast to infoROM object version which is only an indicator of supported features. Version string will not exceed 16 characters in length (including the `NULL` terminator). See [nvmConstants::NVML\\_DEVICE\\_INFOROM\\_VERSION\\_BUFFER\\_SIZE](#).

### See also:

[nvmDeviceGetInforomVersion](#)

## `nvmlReturn_t nvmlDeviceGetInforomConfigurationChecksum (nvmlDevice_t device, unsigned int *checksum)`

### Parameters

#### **device**

The identifier of the target device

#### **checksum**

Reference in which to return the infoROM configuration checksum

### Returns

- ▶ `NVML_SUCCESS` if checksum has been set
- ▶ `NVML_ERROR_CORRUPTED_INFOROM` if the device's checksum couldn't be retrieved due to infoROM corruption
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if checksum is NULL
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the checksum of the configuration stored in the device's infoROM.

For all products with an inforom.

Can be used to make sure that two GPUs have the exact same configuration. Current checksum takes into account configuration stored in PWR and ECC infoROM objects. Checksum can change between driver releases or when user changes configuration (e.g. disable/enable ECC)

## `nvmlReturn_t nvmlDeviceValidateInforom (nvmlDevice_t device)`

### Parameters

#### **device**

The identifier of the target device

**Returns**

- ▶ NVML\_SUCCESS if infoROM is not corrupted
- ▶ NVML\_ERROR\_CORRUPTED\_INFOROM if the device's infoROM is corrupted
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Reads the infoROM from the flash and verifies the checksums.

For all products with an inforom.

**`nvmlReturn_t nvmlDeviceGetLastBBXFlushTime`**  
**`(nvmlDevice_t device, unsigned long long *timestamp,`**  
**`unsigned long *durationUs)`**

**Parameters****device**

The identifier of the target device

**timestamp**

The start timestamp of the last BBX Flush

**durationUs**

The duration (us) of the last BBX Flush

**Returns**

- ▶ NVML\_SUCCESS if timestamp and durationUs are successfully retrieved
- ▶ NVML\_ERROR\_NOT\_READY if the BBX object has not been flushed yet
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not have an infoROM
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the timestamp and the duration of the last flush of the BBX (blackbox) infoROM object during the current run.

For all products with an inforom.

**See also:**

[nvmlDeviceGetInforomVersion](#)

## **nvmlReturn\_t nvmlDeviceGetDisplayMode (nvmlDevice\_t device, nvmlEnableState\_t \*display)**

**Parameters****device**

The identifier of the target device

**display**

Reference in which to return the display mode

**Returns**

- ▶ NVML\_SUCCESS if display has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or display is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the display mode for the device.

For all products.

This method indicates whether a physical display (e.g. monitor) is currently connected to any of the device's connectors.

See [nvmlEnableState\\_t](#) for details on allowed modes.

## **nvmlReturn\_t nvmlDeviceGetDisplayActive (nvmlDevice\_t device, nvmlEnableState\_t \*isActive)**

**Parameters****device**

The identifier of the target device

**isActive**

Reference in which to return the display active state

**Returns**

- ▶ NVML\_SUCCESS if isActive has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or isActive is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the display active state for the device.

For all products.

This method indicates whether a display is initialized on the device. For example whether X Server is attached to this device and has allocated memory for the screen.

Display can be active even when no monitor is physically attached.

See [nvmlEnableState\\_t](#) for details on allowed modes.

## nvmlReturn\_t nvmlDeviceGetPersistenceMode (nvmlDevice\_t device, nvmlEnableState\_t \*mode)

**Parameters****device**

The identifier of the target device

**mode**

Reference in which to return the current driver persistence mode

**Returns**

- ▶ NVML\_SUCCESS if mode has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or mode is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the persistence mode associated with this device.

For all products. For Linux only.

When driver persistence mode is enabled the driver software state is not torn down when the last client disconnects. By default this feature is disabled.

See [nvmlEnableState\\_t](#) for details on allowed modes.

**See also:**

[nvmlDeviceSetPersistenceMode\(\)](#)

## **nvmlReturn\_t nvmlDeviceGetPciInfoExt (nvmlDevice\_t device, nvmlPciInfoExt\_t \*pci)**

### **Parameters**

#### **device**

The identifier of the target device

#### **pci**

Reference in which to return the PCI info

### **Returns**

- ▶ NVML\_SUCCESS if pci has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or pci is NULL
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### **Description**

Retrieves PCI attributes of this device.

For all products.

See [nvmlPciInfoExt\\_t](#) for details on the available PCI info.

## **nvmlReturn\_t nvmlDeviceGetPciInfo\_v3 (nvmlDevice\_t device, nvmlPciInfo\_t \*pci)**

### **Parameters**

#### **device**

The identifier of the target device



**pci**

Reference in which to return the PCI info

**Returns**

- ▶ NVML\_SUCCESS if pci has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or pci is NULL
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the PCI attributes of this device.

For all products.

See [nvmlPciInfo\\_t](#) for details on the available PCI info.

## nvmlReturn\_t nvmlDeviceGetMaxPcieLinkGeneration (nvmlDevice\_t device, unsigned int \*maxLinkGen)

**Parameters****device**

The identifier of the target device

**maxLinkGen**

Reference in which to return the max PCIe link generation

**Returns**

- ▶ NVML\_SUCCESS if maxLinkGen has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or maxLinkGen is null
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if PCIe link information is not available
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the maximum PCIe link generation possible with this device and system

I.E. for a generation 2 PCIe device attached to a generation 1 PCIe bus the max link generation this function will report is generation 1.

For Fermi or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetGpuMaxPcieLinkGeneration` (`nvmlDevice_t device`, `unsigned int *maxLinkGenDevice`)

### Parameters

#### **device**

The identifier of the target device

#### **maxLinkGenDevice**

Reference in which to return the max PCIe link generation

### Returns

- ▶ `NVML_SUCCESS` if `maxLinkGenDevice` has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or `maxLinkGenDevice` is null
- ▶ `NVML_ERROR_NOT_SUPPORTED` if PCIe link information is not available
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the maximum PCIe link generation supported by this device

For Fermi or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetMaxPcieLinkWidth` (`nvmlDevice_t device`, `unsigned int *maxLinkWidth`)

### Parameters

#### **device**

The identifier of the target device

#### **maxLinkWidth**

Reference in which to return the max PCIe link generation

### Returns

- ▶ `NVML_SUCCESS` if `maxLinkWidth` has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized

- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or `maxLinkWidth` is null
- ▶ `NVML_ERROR_NOT_SUPPORTED` if PCIe link information is not available
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the maximum PCIe link width possible with this device and system

I.E. for a device with a 16x PCIe bus width attached to a 8x PCIe system bus this function will report a max link width of 8.

For Fermi or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetCurrPcieLinkGeneration` (`nvmlDevice_t device`, `unsigned int *currLinkGen`)

### Parameters

#### `device`

The identifier of the target device

#### `currLinkGen`

Reference in which to return the current PCIe link generation

### Returns

- ▶ `NVML_SUCCESS` if `currLinkGen` has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or `currLinkGen` is null
- ▶ `NVML_ERROR_NOT_SUPPORTED` if PCIe link information is not available
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the current PCIe link generation

For Fermi or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetCurrPcieLinkWidth` (`nvmlDevice_t device`, `unsigned int *currLinkWidth`)

### Parameters

#### `device`

The identifier of the target device

#### `currLinkWidth`

Reference in which to return the current PCIe link generation

### Returns

- ▶ `NVML_SUCCESS` if `currLinkWidth` has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid or `currLinkWidth` is null
- ▶ `NVML_ERROR_NOT_SUPPORTED` if PCIe link information is not available
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the current PCIe link width

For Fermi or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetPcieThroughput` (`nvmlDevice_t device`, `nvmlPcieUtilCounter_t counter`, `unsigned int *value`)

### Parameters

#### `device`

The identifier of the target device

#### `counter`

The specific counter that should be queried `nvmlPcieUtilCounter_t`

#### `value`

Reference in which to return throughput in KB/s

### Returns

- ▶ `NVML_SUCCESS` if `value` has been set

- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device or counter is invalid, or value is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve PCIe utilization information. This function is querying a byte counter over a 20ms interval and thus is the PCIe throughput over that interval.

For Maxwell or newer fully supported devices.

This method is not supported in virtual machines running virtual GPU (vGPU).

## `nvmlReturn_t nvmlDeviceGetPcieReplayCounter` (`nvmlDevice_t device`, `unsigned int *value`)

### Parameters

#### `device`

The identifier of the target device

#### `value`

Reference in which to return the counter's value

### Returns

- ▶ `NVML_SUCCESS` if value has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, or value is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the PCIe replay counter.

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetClockInfo (nvmlDevice_t device, nvmlClockType_t type, unsigned int *clock)`

### Parameters

**device**

The identifier of the target device

**type**

Identify which clock domain to query

**clock**

Reference in which to return the clock speed in MHz

### Returns

- ▶ `NVML_SUCCESS` if clock has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or clock is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device cannot report the specified clock
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the current clock speeds for the device.

For Fermi or newer fully supported devices.

See `nvmlClockType_t` for details on available clock information.

## `nvmlReturn_t nvmlDeviceGetMaxClockInfo (nvmlDevice_t device, nvmlClockType_t type, unsigned int *clock)`

### Parameters

**device**

The identifier of the target device

**type**

Identify which clock domain to query

**clock**

Reference in which to return the clock speed in MHz

**Returns**

- ▶ NVML\_SUCCESS if clock has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or clock is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device cannot report the specified clock
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the maximum clock speeds for the device.

For Fermi or newer fully supported devices.

See [nvmlClockType\\_t](#) for details on available clock information.



On GPUs from Fermi family current P0 clocks (reported by [nvmlDeviceGetClockInfo](#)) can differ from max clocks by few MHz.

## **`nvmlReturn_t nvmlDeviceGetGpcClkVfOffset`** **(`nvmlDevice_t device, int *offset`)**

**Parameters****device**

The identifier of the target device

**offset**

The retrieved GPCCLK VF offset value

**Returns**

- ▶ NVML\_SUCCESS if offset has been successfully queried
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or offset is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve the GPCCLK VF offset value

**nvmlReturn\_t nvmlDeviceGetApplicationsClock**  
 (nvmlDevice\_t device, nvmlClockType\_t clockType,  
 unsigned int \*clockMHz)

#### Parameters

##### device

The identifier of the target device

##### clockType

Identify which clock domain to query

##### clockMHz

Reference in which to return the clock in MHz

#### Returns

- ▶ NVML\_SUCCESS if clockMHz has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or clockMHz is NULL or clockType is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

#### Description

Retrieves the current setting of a clock that applications will use unless an overspec situation occurs. Can be changed using [nvmlDeviceSetApplicationsClocks](#).

For Kepler or newer fully supported devices.

**nvmlReturn\_t nvmlDeviceGetDefaultApplicationsClock**  
 (nvmlDevice\_t device, nvmlClockType\_t clockType,  
 unsigned int \*clockMHz)

#### Parameters

##### device

The identifier of the target device

##### clockType

Identify which clock domain to query



**clockMHz**

Reference in which to return the default clock in MHz

**Returns**

- ▶ NVML\_SUCCESS if clockMHz has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or clockMHz is NULL or clockType is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the default applications clock that GPU boots with or defaults to after [nvmlDeviceResetApplicationsClocks](#) call.

For Kepler or newer fully supported devices.

**See also:**

[nvmlDeviceGetApplicationsClock](#)

**[nvmlReturn\\_t nvmlDeviceGetClock \(nvmlDevice\\_t device, nvmlClockType\\_t clockType, nvmlClockId\\_t clockId, unsigned int \\*clockMHz\)](#)**

**Parameters****device**

The identifier of the target device

**clockType**

Identify which clock domain to query

**clockId**

Identify which clock in the domain to query

**clockMHz**

Reference in which to return the clock in MHz

**Returns**

- ▶ NVML\_SUCCESS if clockMHz has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized

- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or clockMHz is NULL or clockType is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the clock speed for the clock specified by the clock type and clock ID.

For Kepler or newer fully supported devices.

```
nvmlReturn_t nvmlDeviceGetMaxCustomerBoostClock  
(nvmlDevice_t device, nvmlClockType_t clockType,  
unsigned int *clockMHz)
```

### Parameters

#### **device**

The identifier of the target device

#### **clockType**

Identify which clock domain to query

#### **clockMHz**

Reference in which to return the clock in MHz

### Returns

- ▶ NVML\_SUCCESS if clockMHz has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or clockMHz is NULL or clockType is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device or the clockType on this device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the customer defined maximum boost clock speed specified by the given clock type.

For Pascal or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetSupportedMemoryClocks` (`nvmlDevice_t device`, `unsigned int *count`, `unsigned int *clocksMHz`)

### Parameters

#### `device`

The identifier of the target device

#### `count`

Reference in which to provide the `clocksMHz` array size, and to return the number of elements

#### `clocksMHz`

Reference in which to return the clock in MHz

### Returns

- ▶ `NVML_SUCCESS` if `count` and `clocksMHz` have been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or `count` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `count` is too small (`count` is set to the number of required elements)
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the list of possible memory clocks that can be used as an argument for `nvmlDeviceSetApplicationsClocks`.

For Kepler or newer fully supported devices.

### See also:

[nvmlDeviceSetApplicationsClocks](#)

[nvmlDeviceGetSupportedGraphicsClocks](#)

## `nvmlReturn_t nvmlDeviceGetSupportedGraphicsClocks` (`nvmlDevice_t device`, `unsigned int memoryClockMHz`, `unsigned int *count`, `unsigned int *clocksMHz`)

### Parameters

#### **device**

The identifier of the target device

#### **memoryClockMHz**

Memory clock for which to return possible graphics clocks

#### **count**

Reference in which to provide the `clocksMHz` array size, and to return the number of elements

#### **clocksMHz**

Reference in which to return the clocks in MHz

### Returns

- ▶ `NVML_SUCCESS` if `count` and `clocksMHz` have been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_NOT_FOUND` if the specified `memoryClockMHz` is not a supported frequency
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid or `clock` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `count` is too small
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the list of possible graphics clocks that can be used as an argument for [nvmlDeviceSetApplicationsClocks](#).

For Kepler or newer fully supported devices.

### See also:

[nvmlDeviceSetApplicationsClocks](#)

[nvmlDeviceGetSupportedMemoryClocks](#)

## `nvmlReturn_t nvmlDeviceGetAutoBoostedClocksEnabled` (`nvmlDevice_t device`, `nvmlEnableState_t *isEnabled`, `nvmlEnableState_t *defaultIsEnabled`)

### Parameters

#### **device**

The identifier of the target device

#### **isEnabled**

Where to store the current state of Auto Boosted clocks of the target device

#### **defaultIsEnabled**

Where to store the default Auto Boosted clocks behavior of the target device that the device will revert to when no applications are using the GPU

### Returns

- ▶ `NVML_SUCCESS` If `isEnabled` has been set with the Auto Boosted clocks state of device
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid or `isEnabled` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support Auto Boosted clocks
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the current state of Auto Boosted clocks on a device and store it in `isEnabled`

For Kepler or newer fully supported devices.

Auto Boosted clocks are enabled by default on some hardware, allowing the GPU to run at higher clock rates to maximize performance as thermal limits allow.

On Pascal and newer hardware, Auto Boosted clocks are controlled through application clocks. Use `nvmlDeviceSetApplicationsClocks` and `nvmlDeviceResetApplicationsClocks` to control Auto Boost behavior.

## `nvmlReturn_t nvmlDeviceGetFanSpeed (nvmlDevice_t device, unsigned int *speed)`

### Parameters

#### **device**

The identifier of the target device

#### **speed**

Reference in which to return the fan speed percentage

### Returns

- ▶ `NVML_SUCCESS` if speed has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or speed is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not have a fan
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the intended operating speed of the device's fan.

Note: The reported speed is the intended fan speed. If the fan is physically blocked and unable to spin, the output will not match the actual fan speed.

For all discrete products with dedicated fans.

The fan speed is expressed as a percentage of the product's maximum noise tolerance fan speed. This value may exceed 100% in certain cases.

## `nvmlReturn_t nvmlDeviceGetFanSpeed_v2 (nvmlDevice_t device, unsigned int fan, unsigned int *speed)`

### Parameters

#### **device**

The identifier of the target device

#### **fan**

The index of the target fan, zero indexed.

**speed**

Reference in which to return the fan speed percentage

**Returns**

- ▶ NVML\_SUCCESS if speed has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, fan is not an acceptable index, or speed is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not have a fan or is newer than Maxwell
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the intended operating speed of the device's specified fan.

Note: The reported speed is the intended fan speed. If the fan is physically blocked and unable to spin, the output will not match the actual fan speed.

For all discrete products with dedicated fans.

The fan speed is expressed as a percentage of the product's maximum noise tolerance fan speed. This value may exceed 100% in certain cases.

## nvmlReturn\_t nvmlDeviceGetTargetFanSpeed (nvmlDevice\_t device, unsigned int fan, unsigned int \*targetSpeed)

**Parameters****device**

The identifier of the target device

**fan**

The index of the target fan, zero indexed.

**targetSpeed**

Reference in which to return the fan speed percentage

**Returns**

- ▶ NVML\_SUCCESS if speed has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized

- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, fan is not an acceptable index, or speed is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not have a fan or is newer than Maxwell
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the intended target speed of the device's specified fan.

Normally, the driver dynamically adjusts the fan based on the needs of the GPU. But when user set fan speed using `nvmlDeviceSetFanSpeed_v2`, the driver will attempt to make the fan achieve the setting in `nvmlDeviceSetFanSpeed_v2`. The actual current speed of the fan is reported in `nvmlDeviceGetFanSpeed_v2`.

For all discrete products with dedicated fans.

The fan speed is expressed as a percentage of the product's maximum noise tolerance fan speed. This value may exceed 100% in certain cases.

## `nvmlReturn_t nvmlDeviceGetMinMaxFanSpeed` (`nvmlDevice_t device`, `unsigned int *minSpeed`, `unsigned int *maxSpeed`)

### Parameters

#### **device**

The identifier of the target device

#### **minSpeed**

The minimum speed allowed to set

#### **maxSpeed**

The maximum speed allowed to set

### Description

Retrieves the min and max fan speed that user can set for the GPU fan.

For all cuda-capable discrete products with fans

return `NVML_SUCCESS` if speed has been adjusted

`NVML_ERROR_UNINITIALIZED` if the library has not been successfully

initialized `NVML_ERROR_INVALID_ARGUMENT` if device is invalid

`NVML_ERROR_NOT_SUPPORTED` if the device does not support this (doesn't have fans) `NVML_ERROR_UNKNOWN` on any unexpected error



## `nvmlReturn_t nvmlDeviceGetFanControlPolicy_v2` (`nvmlDevice_t device`, `unsigned int fan`, `nvmlFanControlPolicy_t *policy`)

### Description

Gets current fan control policy.

For Maxwell or newer fully supported devices.

For all cuda-capable discrete products with fans

`device` The identifier of the target device  
`policy` Reference in which to return the fan control policy

return `NVML_SUCCESS` if `policy` has been populated

`NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized

`NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid or `policy` is null or the fan given doesn't reference a fan that exists. `NVML_ERROR_NOT_SUPPORTED` if the device is older than Maxwell `NVML_ERROR_UNKNOWN` on any unexpected error

## `nvmlReturn_t nvmlDeviceGetNumFans` (`nvmlDevice_t device`, `unsigned int *numFans`)

### Parameters

#### `device`

The identifier of the target device

#### `numFans`

The number of fans

### Returns

- ▶ `NVML_SUCCESS` if fan number query was successful
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid or `numFans` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not have a fan
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the number of fans on the device.

For all discrete products with dedicated fans.

```
nvmlReturn_t nvmlDeviceGetTemperature  
(nvmlDevice_t device, nvmlTemperatureSensors_t  
sensorType, unsigned int *temp)
```

### Parameters

#### **device**

The identifier of the target device

#### **sensorType**

Flag that indicates which sensor reading to retrieve

#### **temp**

Reference in which to return the temperature reading

### Returns

- ▶ NVML\_SUCCESS if temp has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, sensorType is invalid or temp is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not have the specified sensor
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the current temperature readings for the device, in degrees C.

For all products.

See [nvmlTemperatureSensors\\_t](#) for details on available temperature sensors.

```
nvmlReturn_t nvmlDeviceGetTemperatureThreshold  
(nvmlDevice_t device, nvmlTemperatureThresholds_t  
thresholdType, unsigned int *temp)
```

### Parameters

#### **device**

The identifier of the target device

**thresholdType**

The type of threshold value queried

**temp**

Reference in which to return the temperature reading

**Returns**

- ▶ NVML\_SUCCESS if temp has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, thresholdType is invalid or temp is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not have a temperature sensor or is unsupported
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the temperature threshold for the GPU with the specified threshold type in degrees C.

For Kepler or newer fully supported devices.

See [nvmlTemperatureThresholds\\_t](#) for details on available temperature thresholds.

Note: This API is no longer the preferred interface for retrieving the following temperature thresholds on Ada and later architectures:

NVML\_TEMPERATURE\_THRESHOLD\_SHUTDOWN,  
 NVML\_TEMPERATURE\_THRESHOLD\_SLOWDOWN,  
 NVML\_TEMPERATURE\_THRESHOLD\_MEM\_MAX and  
 NVML\_TEMPERATURE\_THRESHOLD\_GPU\_MAX.

Support for reading these temperature thresholds for Ada and later architectures would be removed from this API in future releases. Please use [nvmlDeviceGetFieldValues](#) with NVML\_FI\_DEV\_TEMPERATURE\_\* fields to retrieve temperature thresholds on these architectures.

## `nvmlReturn_t nvmlDeviceGetThermalSettings` (`nvmlDevice_t device`, `unsigned int sensorIndex`, `nvmlGpuThermalSettings_t *pThermalSettings`)

### Parameters

#### `device`

The identifier of the target device

#### `sensorIndex`

The index of the thermal sensor

#### `pThermalSettings`

Reference in which to return the thermal sensor information

### Returns

- ▶ `NVML_SUCCESS` if `pThermalSettings` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid or `pThermalSettings` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Used to execute a list of thermal system instructions.

## `nvmlReturn_t nvmlDeviceGetPerformanceState` (`nvmlDevice_t device`, `nvmlPstates_t *pState`)

### Parameters

#### `device`

The identifier of the target device

#### `pState`

Reference in which to return the performance state reading

### Returns

- ▶ `NVML_SUCCESS` if `pState` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized

- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or pState is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the current performance state for the device.

For Fermi or newer fully supported devices.

See [nvmlPstates\\_t](#) for details on allowed performance states.

## **`nvmlReturn_t nvmlDeviceGetCurrentClocksEventReasons`** **(`nvmlDevice_t` device, unsigned long long** **\*`clocksEventReasons`)**

### Parameters

#### **device**

The identifier of the target device

#### **clocksEventReasons**

Reference in which to return bitmask of active clocks event reasons

### Returns

- ▶ NVML\_SUCCESS if clocksEventReasons has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or clocksEventReasons is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves current clocks event reasons.

For all fully supported products.



More than one bit can be enabled at the same time. Multiple reasons can be affecting clocks at once.

**See also:**[NvmlClocksEventReasons](#)[nvmlDeviceGetSupportedClocksEventReasons](#)

**nvmlReturn\_t**  
**nvmlDeviceGetCurrentClocksThrottleReasons**  
 (nvmlDevice\_t device, unsigned long long  
 \*clocksThrottleReasons)

**Description**

Deprecated Use [nvmlDeviceGetCurrentClocksEventReasons](#) instead

**nvmlReturn\_t**  
**nvmlDeviceGetSupportedClocksEventReasons**  
 (nvmlDevice\_t device, unsigned long long  
 \*supportedClocksEventReasons)

**Parameters****device**

The identifier of the target device

**supportedClocksEventReasons**

Reference in which to return bitmask of supported clocks event reasons

**Returns**

- ▶ NVML\_SUCCESS if supportedClocksEventReasons has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or supportedClocksEventReasons is NULL
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves bitmask of supported clocks event reasons that can be returned by [nvmlDeviceGetCurrentClocksEventReasons](#)

For all fully supported products.

This method is not supported in virtual machines running virtual GPU (vGPU).

**See also:**

[NvmlClocksEventReasons](#)

[nvmlDeviceGetCurrentClocksEventReasons](#)

**nvmlReturn\_t  
nvmlDeviceGetSupportedClocksThrottleReasons  
(nvmlDevice\_t device, unsigned long long  
\*supportedClocksThrottleReasons)**

**Description**

Deprecated Use [nvmlDeviceGetSupportedClocksEventReasons](#) instead

**nvmlReturn\_t nvmlDeviceGetPowerState (nvmlDevice\_t  
device, nvmlPstates\_t \*pState)**

**Parameters**

**device**

The identifier of the target device

**pState**

Reference in which to return the performance state reading

**Returns**

- ▶ NVML\_SUCCESS if pState has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or pState is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Deprecated: Use [nvmlDeviceGetPerformanceState](#). This function exposes an incorrect generalization.

Retrieve the current performance state for the device.

For Fermi or newer fully supported devices.

See `nvmlPstates_t` for details on allowed performance states.

## `nvmlReturn_t nvmlDeviceGetDynamicPstatesInfo` (`nvmlDevice_t device`, `nvmlGpuDynamicPstatesInfo_t *pDynamicPstatesInfo`)

### Parameters

`device`

`pDynamicPstatesInfo`

### Returns

- ▶ `NVML_SUCCESS` if `pDynamicPstatesInfo` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid or `pDynamicPstatesInfo` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve performance monitor samples from the associated subdevice.

## `nvmlReturn_t nvmlDeviceGetMemClkVfOffset` (`nvmlDevice_t device`, `int *offset`)

### Parameters

`device`

The identifier of the target device

`offset`

The retrieved MemClk VF offset value

### Returns

- ▶ `NVML_SUCCESS` if `offset` has been successfully queried
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid or `offset` is `NULL`



- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieve the MemClk (Memory Clock) VF offset value.

```
nvmlReturn_t nvmlDeviceGetMinMaxClockOfPState
(nvmlDevice_t device, nvmlClockType_t type,
nvmlPstates_t pstate, unsigned int *minClockMHz,
unsigned int *maxClockMHz)
```

### Parameters

#### **device**

The identifier of the target device

#### **type**

Clock domain

#### **pstate**

PState to query

#### **minClockMHz**

Reference in which to return min clock frequency

#### **maxClockMHz**

Reference in which to return max clock frequency

### Returns

- ▶ NVML\_SUCCESS if everything worked
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, type or pstate are invalid or both minClockMHz and maxClockMHz are NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature

### Description

Retrieve min and max clocks of some clock domain for a given PState

```
nvmlReturn_t
nvmlDeviceGetSupportedPerformanceStates
```

## `(nvmlDevice_t device, nvmlPstates_t *pstates, unsigned int size)`

### Parameters

#### **device**

The identifier of the target device

#### **pstates**

Container to return the list of performance states supported by device

#### **size**

Size of the supplied pstates array in bytes

### Returns

- ▶ NVML\_SUCCESS if pstates array has been retrieved
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if the the container supplied was not large enough to hold the resulting list
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device or pstates is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support performance state readings
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Get all supported Performance States (P-States) for the device.

The returned array would contain a contiguous list of valid P-States supported by the device. If the number of supported P-States is fewer than the size of the array supplied missing elements would contain NVML\_PSTATE\_UNKNOWN.

The number of elements in the returned list will never exceed NVML\_MAX\_GPU\_PERF\_PSTATES.

## `nvmlReturn_t nvmlDeviceGetGpcClkMinMaxVfOffset` `(nvmlDevice_t device, int *minOffset, int *maxOffset)`

### Parameters

#### **device**

The identifier of the target device

#### **minOffset**

The retrieved GPCCLK VF min offset value

**maxOffset**

The retrieved GPCCLK VF max offset value

**Returns**

- ▶ NVML\_SUCCESS if offset has been successfully queried
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or offset is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve the GPCCLK min max VF offset value.

**`nvmlReturn_t nvmlDeviceGetMemClkMinMaxVfOffset  
(nvmlDevice_t device, int *minOffset, int *maxOffset)`**

**Parameters****device**

The identifier of the target device

**minOffset**

The retrieved MemClk VF min offset value

**maxOffset**

The retrieved MemClk VF max offset value

**Returns**

- ▶ NVML\_SUCCESS if offset has been successfully queried
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or offset is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve the MemClk (Memory Clock) min max VF offset value.

## `nvmlReturn_t nvmlDeviceGetPowerManagementMode` (`nvmlDevice_t device, nvmlEnableState_t *mode`)

### Parameters

#### **device**

The identifier of the target device

#### **mode**

Reference in which to return the current power management mode

### Returns

- ▶ `NVML_SUCCESS` if mode has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or mode is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

This API has been deprecated.

Retrieves the power management mode associated with this device.

For products from the Fermi family.

- ▶ Requires `NVML_INFOROM_POWER` version 3.0 or higher.

For from the Kepler or newer families.

- ▶ Does not require `NVML_INFOROM_POWER` object.

This flag indicates whether any power management algorithm is currently active on the device. An enabled state does not necessarily mean the device is being actively throttled -- only that that the driver will do so if the appropriate conditions are met.

See [nvmlEnableState\\_t](#) for details on allowed modes.

## `nvmlReturn_t nvmlDeviceGetPowerManagementLimit` (`nvmlDevice_t device, unsigned int *limit`)

### Parameters

**device**

The identifier of the target device

**limit**

Reference in which to return the power management limit in milliwatts

### Returns

- ▶ `NVML_SUCCESS` if limit has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or limit is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the power management limit associated with this device.

For Fermi or newer fully supported devices.

The power limit defines the upper boundary for the card's power draw. If the card's total power draw reaches this limit the power management algorithm kicks in.

This reading is only available if power management mode is supported. See [nvmlDeviceGetPowerManagementMode](#).

## `nvmlReturn_t` `nvmlDeviceGetPowerManagementLimitConstraints` (`nvmlDevice_t device, unsigned int *minLimit, unsigned int *maxLimit`)

### Parameters

**device**

The identifier of the target device

**minLimit**

Reference in which to return the minimum power management limit in milliwatts

**maxLimit**

Reference in which to return the maximum power management limit in milliwatts

**Returns**

- ▶ NVML\_SUCCESS if minLimit and maxLimit have been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or minLimit or maxLimit is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves information about possible values of power management limits on this device. For Kepler or newer fully supported devices.

**See also:**

[nvmlDeviceSetPowerManagementLimit](#)

## nvmlReturn\_t nvmlDeviceGetPowerManagementDefaultLimit (nvmlDevice\_t device, unsigned int \*defaultLimit)

**Parameters****device**

The identifier of the target device

**defaultLimit**

Reference in which to return the default power management limit in milliwatts

**Returns**

- ▶ NVML\_SUCCESS if defaultLimit has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or defaultLimit is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible

- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves default power management limit on this device, in milliwatts. Default power management limit is a power management limit that the device boots with.

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetPowerUsage (nvmlDevice_t device, unsigned int *power)`

### Parameters

#### **device**

The identifier of the target device

#### **power**

Reference in which to return the power usage information

### Returns

- ▶ NVML\_SUCCESS if power has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or power is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support power readings
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves power usage for this GPU in milliwatts and its associated circuitry (e.g. memory)

For Fermi or newer fully supported devices.

On Fermi and Kepler GPUs the reading is accurate to within +/- 5% of current power draw. On Ampere (except GA100) or newer GPUs, the API returns power averaged over 1 sec interval. On GA100 and older architectures, instantaneous power is returned.

See [NVML\\_FI\\_DEV\\_POWER\\_AVERAGE](#) and [NVML\\_FI\\_DEV\\_POWER\\_INSTANT](#) to query specific power values.

It is only available if power management mode is supported. See [nvmlDeviceGetPowerManagementMode](#).

## `nvmlReturn_t nvmlDeviceGetTotalEnergyConsumption` (`nvmlDevice_t device`, `unsigned long long *energy`)

### Parameters

#### **device**

The identifier of the target device

#### **energy**

Reference in which to return the energy consumption information

### Returns

- ▶ `NVML_SUCCESS` if energy has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or energy is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support energy readings
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves total energy consumption for this GPU in millijoules (mJ) since the driver was last reloaded

For Volta or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetEnforcedPowerLimit` (`nvmlDevice_t device`, `unsigned int *limit`)

### Parameters

#### **device**

The device to communicate with

#### **limit**

Reference in which to return the power management limit in milliwatts

### Returns

- ▶ `NVML_SUCCESS` if limit has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or limit is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature



- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Get the effective power limit that the driver enforces after taking into account all limiters

Note: This can be different from the `nvmlDeviceGetPowerManagementLimit` if other limits are set elsewhere This includes the out of band power limit interface

For Kepler or newer fully supported devices.

```
nvmlReturn_t nvmlDeviceGetGpuOperationMode
(nvmlDevice_t device, nvmlGpuOperationMode_t
*current, nvmlGpuOperationMode_t *pending)
```

### Parameters

#### **device**

The identifier of the target device

#### **current**

Reference in which to return the current GOM

#### **pending**

Reference in which to return the pending GOM

### Returns

- ▶ `NVML_SUCCESS` if mode has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or current or pending is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the current GOM and pending GOM (the one that GPU will switch to after reboot).

For GK110 M-class and X-class Tesla products from the Kepler family. Modes `NVML_GOM_LOW_DP` and `NVML_GOM_ALL_ON` are supported on fully supported GeForce products. Not supported on Quadro and Tesla C-class products.

**See also:**[nvmlGpuOperationMode\\_t](#)[nvmlDeviceSetGpuOperationMode](#)

## **`nvmlReturn_t nvmlDeviceGetMemoryInfo (nvmlDevice_t device, nvmlMemory_t *memory)`**

**Parameters****device**

The identifier of the target device

**memory**

Reference in which to return the memory information

**Returns**

- ▶ `NVML_SUCCESS` if memory has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or memory is `NULL`
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

**Description**

Retrieves the amount of used, free, reserved and total memory available on the device, in bytes. The reserved amount is supported on version 2 only.

For all products.

Enabling ECC reduces the amount of total available memory, due to the extra required parity bits. Under WDDM most device memory is allocated and managed on startup by Windows.

Under Linux and Windows TCC, the reported amount of used memory is equal to the sum of memory allocated by all active channels on the device.

See [nvmlMemory\\_v2\\_t](#) for details on available memory info.



- ▶ In MIG mode, if device handle is provided, the API returns aggregate information, only if the caller has appropriate privileges. Per-instance information can be queried by using specific MIG device handles.

- ▶ `nvmlDeviceGetMemoryInfo_v2` adds additional memory information.
- ▶ On systems where GPUs are NUMA nodes, the accuracy of FB memory utilization provided by this API depends on the memory accounting of the operating system. This is because FB memory is managed by the operating system instead of the NVIDIA GPU driver. Typically, pages allocated from FB memory are not released even after the process terminates to enhance performance. In scenarios where the operating system is under memory pressure, it may resort to utilizing FB memory. Such actions can result in discrepancies in the accuracy of memory reporting.

## `nvmlReturn_t nvmlDeviceGetComputeMode` (`nvmlDevice_t device`, `nvmlComputeMode_t *mode`)

### Parameters

#### **device**

The identifier of the target device

#### **mode**

Reference in which to return the current compute mode

### Returns

- ▶ `NVML_SUCCESS` if mode has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or mode is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the current compute mode for the device.

For all products.

See `nvmlComputeMode_t` for details on allowed compute modes.

### See also:

`nvmlDeviceSetComputeMode()`

## `nvmlReturn_t nvmlDeviceGetCudaComputeCapability (nvmlDevice_t device, int *major, int *minor)`

### Parameters

#### **device**

The identifier of the target device

#### **major**

Reference in which to return the major CUDA compute capability

#### **minor**

Reference in which to return the minor CUDA compute capability

### Returns

- ▶ `NVML_SUCCESS` if major and minor have been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or major or minor are `NULL`
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the CUDA compute capability of the device.

For all products.

Returns the major and minor compute capability version numbers of the device. The major and minor versions are equivalent to the `CU_DEVICE_ATTRIBUTE_COMPUTE_CAPABILITY_MINOR` and `CU_DEVICE_ATTRIBUTE_COMPUTE_CAPABILITY_MAJOR` attributes that would be returned by CUDA's `cuDeviceGetAttribute()`.

## `nvmlReturn_t nvmlDeviceGetEccMode (nvmlDevice_t device, nvmlEnableState_t *current, nvmlEnableState_t *pending)`

### Parameters

#### **device**

The identifier of the target device

**current**

Reference in which to return the current ECC mode

**pending**

Reference in which to return the pending ECC mode

**Returns**

- ▶ NVML\_SUCCESS if current and pending have been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or either current or pending is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the current and pending ECC modes for the device.

For Fermi or newer fully supported devices. Only applicable to devices with ECC. Requires NVML\_INFOROM\_ECC version 1.0 or higher.

Changing ECC modes requires a reboot. The "pending" ECC mode refers to the target mode following the next reboot.

See [nvmlEnableState\\_t](#) for details on allowed modes.

**See also:**

[nvmlDeviceSetEccMode\(\)](#)

## **`nvmlReturn_t nvmlDeviceGetDefaultEccMode` `(nvmlDevice_t device, nvmlEnableState_t *defaultMode)`**

**Parameters****device**

The identifier of the target device

**defaultMode**

Reference in which to return the default ECC mode

**Returns**

- ▶ NVML\_SUCCESS if current and pending have been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized

- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or default is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the default ECC modes for the device.

For Fermi or newer fully supported devices. Only applicable to devices with ECC. Requires `NVML_INFOROM_ECC` version 1.0 or higher.

See `nvmlEnableState_t` for details on allowed modes.

### See also:

`nvmlDeviceSetEccMode()`

## `nvmlReturn_t nvmlDeviceGetBoardId (nvmlDevice_t device, unsigned int *boardId)`

### Parameters

#### **device**

The identifier of the target device

#### **boardId**

Reference in which to return the device's board ID

### Returns

- ▶ `NVML_SUCCESS` if `boardId` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or `boardId` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the device `boardId` from 0-N. Devices with the same `boardId` indicate GPUs connected to the same PLX. Use in conjunction with `nvmlDeviceGetMultiGpuBoard()` to decide if they are on the same board as well. The `boardId` returned is a unique ID for the current configuration. Uniqueness and ordering across reboots and system

configurations is not guaranteed (i.e. if a Tesla K40c returns 0x100 and the two GPUs on a Tesla K10 in the same system returns 0x200 it is not guaranteed they will always return those values but they will always be different from each other).

For Fermi or newer fully supported devices.

## **nvmlReturn\_t nvmlDeviceGetMultiGpuBoard (nvmlDevice\_t device, unsigned int \*multiGpuBool)**

### **Parameters**

#### **device**

The identifier of the target device

#### **multiGpuBool**

Reference in which to return a zero or non-zero value to indicate whether the device is on a multi GPU board

### **Returns**

- ▶ NVML\_SUCCESS if multiGpuBool has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or multiGpuBool is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### **Description**

Retrieves whether the device is on a Multi-GPU Board Devices that are on multi-GPU boards will set multiGpuBool to a non-zero value.

For Fermi or newer fully supported devices.

## **nvmlReturn\_t nvmlDeviceGetTotalEccErrors (nvmlDevice\_t device, nvmlMemoryErrorType\_t**

## errorType, nvmlEccCounterType\_t counterType, unsigned long long \*eccCounts)

### Parameters

#### device

The identifier of the target device

#### errorType

Flag that specifies the type of the errors.

#### counterType

Flag that specifies the counter-type of the errors.

#### eccCounts

Reference in which to return the specified ECC errors

### Returns

- ▶ NVML\_SUCCESS if eccCounts has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, errorType or counterType is invalid, or eccCounts is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the total ECC error counts for the device.

For Fermi or newer fully supported devices. Only applicable to devices with ECC. Requires NVML\_INFOROM\_ECC version 1.0 or higher. Requires ECC Mode to be enabled.

The total error count is the sum of errors across each of the separate memory systems, i.e. the total set of errors across the entire device.

See [nvmlMemoryErrorType\\_t](#) for a description of available error types. See [nvmlEccCounterType\\_t](#) for a description of available counter types.

### See also:

[nvmlDeviceClearEccErrorCounts\(\)](#)

## [nvmlReturn\\_t nvmlDeviceGetDetailedEccErrors](#) ([nvmlDevice\\_t](#) device, [nvmlMemoryErrorType\\_t](#)



## errorType, nvmlEccCounterType\_t counterType, nvmlEccErrorCounts\_t \*eccCounts)

### Parameters

#### device

The identifier of the target device

#### errorType

Flag that specifies the type of the errors.

#### counterType

Flag that specifies the counter-type of the errors.

#### eccCounts

Reference in which to return the specified ECC errors

### Returns

- ▶ NVML\_SUCCESS if eccCounts has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, errorType or counterType is invalid, or eccCounts is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the detailed ECC error counts for the device.

**Deprecated** This API supports only a fixed set of ECC error locations. On different GPU architectures different locations are supported. See [nvmlDeviceGetMemoryErrorCounter](#)

For Fermi or newer fully supported devices. Only applicable to devices with ECC. Requires NVML\_INFOROM\_ECC version 2.0 or higher to report aggregate location-based ECC counts. Requires NVML\_INFOROM\_ECC version 1.0 or higher to report all other ECC counts. Requires ECC Mode to be enabled.

Detailed errors provide separate ECC counts for specific parts of the memory system.

Reports zero for unsupported ECC error counters when a subset of ECC error counters are supported.

See [nvmlMemoryErrorType\\_t](#) for a description of available bit types. See [nvmlEccCounterType\\_t](#) for a description of available counter types. See [nvmlEccErrorCounts\\_t](#) for a description of provided detailed ECC counts.

### See also:

```
nvmlDeviceClearEccErrorCounts()
```

```
nvmlReturn_t nvmlDeviceGetMemoryErrorCounter
(nvmlDevice_t device, nvmlMemoryErrorType_t
errorType, nvmlEccCounterType_t counterType,
nvmlMemoryLocation_t locationType, unsigned long long
*count)
```

### Parameters

#### device

The identifier of the target device

#### errorType

Flag that specifies the type of error.

#### counterType

Flag that specifies the counter-type of the errors.

#### locationType

Specifies the location of the counter.

#### count

Reference in which to return the ECC counter

### Returns

- ▶ NVML\_SUCCESS if count has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, errorType, counterType or locationType is invalid, or count is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support ECC error reporting in the specified memory
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the requested memory error counter for the device.

For Fermi or newer fully supported devices. Requires NVML\_INFOROM\_ECC version 2.0 or higher to report aggregate location-based memory error counts. Requires NVML\_INFOROM\_ECC version 1.0 or higher to report all other memory error counts.

Only applicable to devices with ECC.

Requires ECC Mode to be enabled.



On MIG-enabled GPUs, per instance information can be queried using specific MIG device handles. Per instance information is currently only supported for non-DRAM uncorrectable volatile errors. Querying volatile errors using device handles is currently not supported.

See [nvmlMemoryErrorType\\_t](#) for a description of available memory error types. See [nvmlEccCounterType\\_t](#) for a description of available counter types. See [nvmlMemoryLocation\\_t](#) for a description of available counter locations.

## **nvmlReturn\_t nvmlDeviceGetUtilizationRates** (**nvmlDevice\_t device, nvmlUtilization\_t \*utilization**)

### Parameters

#### **device**

The identifier of the target device

#### **utilization**

Reference in which to return the utilization information

### Returns

- ▶ NVML\_SUCCESS if utilization has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or utilization is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the current utilization rates for the device's major subsystems.

For Fermi or newer fully supported devices.

See [nvmlUtilization\\_t](#) for details on available utilization rates.



- ▶ During driver initialization when ECC is enabled one can see high GPU and Memory Utilization readings. This is caused by ECC Memory Scrubbing mechanism that is performed during driver initialization.
- ▶ On MIG-enabled GPUs, querying device utilization rates is not currently supported.

## `nvmlReturn_t nvmlDeviceGetEncoderUtilization` (`nvmlDevice_t device`, `unsigned int *utilization`, `unsigned int *samplingPeriodUs`)

### Parameters

#### **device**

The identifier of the target device

#### **utilization**

Reference to an unsigned int for encoder utilization info

#### **samplingPeriodUs**

Reference to an unsigned int for the sampling period in US

### Returns

- ▶ `NVML_SUCCESS` if utilization has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, utilization is `NULL`, or `samplingPeriodUs` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the current utilization and sampling size in microseconds for the Encoder For Kepler or newer fully supported devices.



On MIG-enabled GPUs, querying encoder utilization is not currently supported.

## `nvmlReturn_t nvmlDeviceGetEncoderCapacity` (`nvmlDevice_t device`, `nvmlEncoderType_t` `encoderQueryType`, `unsigned int *encoderCapacity`)

### Parameters

#### **device**

The identifier of the target device

**encoderQueryType**

Type of encoder to query

**encoderCapacity**

Reference to an unsigned int for the encoder capacity

**Returns**

- ▶ NVML\_SUCCESS if encoderCapacity is fetched
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if encoderCapacity is NULL, or device or encoderQueryType are invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if device does not support the encoder specified in encodeQueryType
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the current capacity of the device's encoder, as a percentage of maximum encoder capacity with valid values in the range 0-100.

For Maxwell or newer fully supported devices.

**nvmlReturn\_t nvmlDeviceGetEncoderStats**  
**(nvmlDevice\_t device, unsigned int \*sessionCount,**  
**unsigned int \*averageFps, unsigned int \*averageLatency)**

**Parameters****device**

The identifier of the target device

**sessionCount**

Reference to an unsigned int for count of active encoder sessions

**averageFps**

Reference to an unsigned int for trailing average FPS of all active sessions

**averageLatency**

Reference to an unsigned int for encode latency in microseconds

**Returns**

- ▶ NVML\_SUCCESS if sessionCount, averageFps and averageLatency is fetched
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized

- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if sessionCount, or device or averageFps, or averageLatency is NULL
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the current encoder statistics for a given device.

For Maxwell or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetEncoderSessions` (`nvmlDevice_t device`, `unsigned int *sessionCount`, `nvmlEncoderSessionInfo_t *sessionInfos`)

### Parameters

#### **device**

The identifier of the target device

#### **sessionCount**

Reference to caller supplied array size, and returns the number of sessions.

#### **sessionInfos**

Reference in which to return the session information

### Returns

- ▶ NVML\_SUCCESS if sessionInfos is fetched
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if sessionCount is too small, array element count is returned in sessionCount
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if sessionCount is NULL.
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by device
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves information about active encoder sessions on a target device.

An array of active encoder sessions is returned in the caller-supplied buffer pointed at by sessionInfos. The array element count is passed in sessionCount, and sessionCount is used to return the number of sessions written to the buffer.

If the supplied buffer is not large enough to accommodate the active session array, the function returns `NVML_ERROR_INSUFFICIENT_SIZE`, with the element count of `nvmlEncoderSessionInfo_t` array required in `sessionCount`. To query the number of active encoder sessions, call this function with `*sessionCount = 0`. The code will return `NVML_SUCCESS` with number of active encoder sessions updated in `*sessionCount`.

For Maxwell or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetDecoderUtilization` (`nvmlDevice_t device`, `unsigned int *utilization`, `unsigned int *samplingPeriodUs`)

### Parameters

#### `device`

The identifier of the target device

#### `utilization`

Reference to an unsigned int for decoder utilization info

#### `samplingPeriodUs`

Reference to an unsigned int for the sampling period in US

### Returns

- ▶ `NVML_SUCCESS` if utilization has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, utilization is `NULL`, or `samplingPeriodUs` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the current utilization and sampling size in microseconds for the Decoder

For Kepler or newer fully supported devices.



On MIG-enabled GPUs, querying decoder utilization is not currently supported.

## `nvmlReturn_t nvmlDeviceGetJpgUtilization` (`nvmlDevice_t device`, `unsigned int *utilization`, `unsigned int *samplingPeriodUs`)

### Parameters

#### **device**

The identifier of the target device

#### **utilization**

Reference to an unsigned int for jpg utilization info

#### **samplingPeriodUs**

Reference to an unsigned int for the sampling period in US

### Returns

- ▶ `NVML_SUCCESS` if utilization has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, utilization is `NULL`, or `samplingPeriodUs` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the current utilization and sampling size in microseconds for the `JPG TURING_OR_NEWER%`



On MIG-enabled GPUs, querying decoder utilization is not currently supported.

## `nvmlReturn_t nvmlDeviceGetOfaUtilization` (`nvmlDevice_t device`, `unsigned int *utilization`, `unsigned int *samplingPeriodUs`)

### Parameters

#### **device**

The identifier of the target device



**utilization**

Reference to an unsigned int for ofa utilization info

**samplingPeriodUs**

Reference to an unsigned int for the sampling period in US

**Returns**

- ▶ NVML\_SUCCESS if utilization has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, utilization is NULL, or samplingPeriodUs is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the current utilization and sampling size in microseconds for the OFA (Optical Flow Accelerator)

TURING\_OR\_NEWER%



On MIG-enabled GPUs, querying decoder utilization is not currently supported.

## **nvmlReturn\_t nvmlDeviceGetFBCStats (nvmlDevice\_t device, nvmlFBCStats\_t \*fbcStats)**

**Parameters****device**

The identifier of the target device

**fbcStats**

Reference to [nvmlFBCStats\\_t](#) structure containing NvFBC stats

**Returns**

- ▶ NVML\_SUCCESS if fbcStats is fetched
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if fbcStats is NULL
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the active frame buffer capture sessions statistics for a given device.

For Maxwell or newer fully supported devices.

```
nvmlReturn_t nvmlDeviceGetFBCSessions
(nvmlDevice_t device, unsigned int *sessionCount,
nvmlFBCSessionInfo_t *sessionInfo)
```

**Parameters****device**

The identifier of the target device

**sessionCount**

Reference to caller supplied array size, and returns the number of sessions.

**sessionInfo**

Reference in which to return the session information

**Returns**

- ▶ NVML\_SUCCESS if sessionInfo is fetched
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if sessionCount is too small, array element count is returned in sessionCount
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if sessionCount is NULL.
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves information about active frame buffer capture sessions on a target device.

An array of active FBC sessions is returned in the caller-supplied buffer pointed at by sessionInfo. The array element count is passed in sessionCount, and sessionCount is used to return the number of sessions written to the buffer.

If the supplied buffer is not large enough to accommodate the active session array, the function returns NVML\_ERROR\_INSUFFICIENT\_SIZE, with the element count of nvmlFBCSessionInfo\_t array required in sessionCount. To query the number of active FBC sessions, call this function with \*sessionCount = 0. The code will return NVML\_SUCCESS with number of active FBC sessions updated in \*sessionCount.

For Maxwell or newer fully supported devices.



hResolution, vResolution, averageFPS and averageLatency data for a FBC session returned in sessionInfo may be zero if there are no new frames captured since the session started.

## `nvmlReturn_t nvmlDeviceGetDriverModel (nvmlDevice_t device, nvmlDriverModel_t *current, nvmlDriverModel_t *pending)`

### Parameters

#### **device**

The identifier of the target device

#### **current**

Reference in which to return the current driver model

#### **pending**

Reference in which to return the pending driver model

### Returns

- ▶ `NVML_SUCCESS` if either current and/or pending have been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or both current and pending are NULL
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the platform is not windows
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the current and pending driver model for the device.

For Fermi or newer fully supported devices. For windows only.

On Windows platforms the device driver can run in either WDDM or WDM (TCC) mode. If a display is attached to the device it must run in WDDM mode. TCC mode is preferred if a display is not attached.

See [nvmlDriverModel\\_t](#) for details on available driver models.

### See also:

[nvmlDeviceSetDriverModel\(\)](#)

## `nvmlReturn_t nvmlDeviceGetVbiosVersion (nvmlDevice_t device, char *version, unsigned int length)`

### Parameters

**device**

The identifier of the target device

**version**

Reference to which to return the VBIOS version

**length**

The maximum allowed length of the string returned in version

### Returns

- ▶ `NVML_SUCCESS` if version has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, or version is `NULL`
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if length is too small
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Get VBIOS version of the device.

For all products.

The VBIOS version may change from time to time. It will not exceed 32 characters in length (including the `NULL` terminator). See [nvmlConstants::NVML\\_DEVICE\\_VBIOS\\_VERSION\\_BUFFER\\_SIZE](#).

## `nvmlReturn_t nvmlDeviceGetBridgeChipInfo (nvmlDevice_t device, nvmlBridgeChipHierarchy_t *bridgeHierarchy)`

### Parameters

**device**

The identifier of the target device

**bridgeHierarchy**

Reference to the returned bridge chip Hierarchy

**Returns**

- ▶ NVML\_SUCCESS if bridge chip exists
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, or bridgeInfo is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if bridge chip not supported on the device
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Get Bridge Chip Information for all the bridge chips on the board.

For all fully supported products. Only applicable to multi-GPU products.

**nvmlReturn\_t**  
**nvmlDeviceGetComputeRunningProcesses\_v3**  
 (nvmlDevice\_t device, unsigned int \*infoCount,  
 nvmlProcessInfo\_t \*infos)

**Parameters****device**

The device handle or MIG device handle

**infoCount**

Reference in which to provide the infos array size, and to return the number of returned elements

**infos**

Reference in which to return the process information

**Returns**

- ▶ NVML\_SUCCESS if infoCount and infos have been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if infoCount indicates that the infos array is too small infoCount will contain minimal amount of space necessary for the call to complete
- ▶ NVML\_ERROR\_NO\_PERMISSION if the user doesn't have permission to perform this operation
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, either of infoCount or infos is NULL

- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by device
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Get information about processes with a compute context on a device

For Fermi or newer fully supported devices.

This function returns information only about compute running processes (e.g. CUDA application which have active context). Any graphics applications (e.g. using OpenGL, DirectX) won't be listed by this function.

To query the current number of running compute processes, call this function with `*infoCount = 0`. The return code will be `NVML_ERROR_INSUFFICIENT_SIZE`, or `NVML_SUCCESS` if none are running. For this call `infos` is allowed to be `NULL`.

The `usedGpuMemory` field returned is all of the memory used by the application.

Keep in mind that information returned by this call is dynamic and the number of elements might change in time. Allocate more space for `infos` table in case new compute processes are spawned.



In MIG mode, if device handle is provided, the API returns aggregate information, only if the caller has appropriate privileges. Per-instance information can be queried by using specific MIG device handles. Querying per-instance information using MIG device handles is not supported if the device is in vGPU Host virtualization mode.

### See also:

[nvmSystemGetProcessName](#)

[nvmReturn\\_t](#)

[nvmDeviceGetGraphicsRunningProcesses\\_v3](#)

([nvmDevice\\_t](#) device, unsigned int \*infoCount, [nvmProcessInfo\\_t](#) \*infos)

### Parameters

#### device

The device handle or MIG device handle

**infoCount**

Reference in which to provide the infos array size, and to return the number of returned elements

**infos**

Reference in which to return the process information

**Returns**

- ▶ NVML\_SUCCESS if infoCount and infos have been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if infoCount indicates that the infos array is too small infoCount will contain minimal amount of space necessary for the call to complete
- ▶ NVML\_ERROR\_NO\_PERMISSION if the user doesn't have permission to perform this operation
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, either of infoCount or infos is NULL
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by device
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Get information about processes with a graphics context on a device

For Kepler or newer fully supported devices.

This function returns information only about graphics based processes (eg. applications using OpenGL, DirectX)

To query the current number of running graphics processes, call this function with \*infoCount = 0. The return code will be NVML\_ERROR\_INSUFFICIENT\_SIZE, or NVML\_SUCCESS if none are running. For this call infos is allowed to be NULL.

The usedGpuMemory field returned is all of the memory used by the application.

Keep in mind that information returned by this call is dynamic and the number of elements might change in time. Allocate more space for infos table in case new graphics processes are spawned.



In MIG mode, if device handle is provided, the API returns aggregate information, only if the caller has appropriate privileges. Per-instance information can be queried by using specific MIG device handles. Querying per-instance information using MIG device handles is not supported if the device is in vGPU Host virtualization mode.

**See also:**

[nvmlSystemGetProcessName](#)

**nvmlReturn\_t****nvmlDeviceGetMPSComputeRunningProcesses\_v3**

(**nvmlDevice\_t** device, unsigned int \*infoCount, **nvmlProcessInfo\_t** \*infos)

**Parameters****device**

The device handle or MIG device handle

**infoCount**

Reference in which to provide the infos array size, and to return the number of returned elements

**infos**

Reference in which to return the process information

**Returns**

- ▶ NVML\_SUCCESS if infoCount and infos have been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if infoCount indicates that the infos array is too small infoCount will contain minimal amount of space necessary for the call to complete
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, either of infoCount or infos is NULL
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by device
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Get information about processes with a MPS compute context on a device

For Volta or newer fully supported devices.

This function returns information only about compute running processes (e.g. CUDA application which have active context) utilizing MPS. Any graphics applications (e.g. using OpenGL, DirectX) won't be listed by this function.



To query the current number of running compute processes, call this function with `*infoCount = 0`. The return code will be `NVML_ERROR_INSUFFICIENT_SIZE`, or `NVML_SUCCESS` if none are running. For this call `infos` is allowed to be `NULL`.

The `usedGpuMemory` field returned is all of the memory used by the application.

Keep in mind that information returned by this call is dynamic and the number of elements might change in time. Allocate more space for `infos` table in case new compute processes are spawned.



In MIG mode, if device handle is provided, the API returns aggregate information, only if the caller has appropriate privileges. Per-instance information can be queried by using specific MIG device handles. Querying per-instance information using MIG device handles is not supported if the device is in vGPU Host virtualization mode.

#### See also:

[nvmlSystemGetProcessName](#)

## `nvmlReturn_t nvmlDeviceGetRunningProcessDetailList` (`nvmlDevice_t device, nvmlProcessDetailList_t *plist`)

### Parameters

#### **device**

The device handle or MIG device handle

#### **plist**

Reference in which to process detail list

### Returns

- ▶ `NVML_SUCCESS` if `plist->numprocArrayEntries` and `plist->procArray` have been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `plist->numprocArrayEntries` indicates that the `plist->procArray` is too small `plist->numprocArrayEntries` will contain minimal amount of space necessary for the call to complete
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, `plist` is `NULL`, `plist->version` is invalid, `plist->mode` is invalid,
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by device

- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Get information about running processes on a device for input context

`HOPPER_OR_NEWER%`

This function returns information only about running processes (e.g. CUDA application which have active context).

To determine the size of the `plist->procArray` array to allocate, call the function with `plist->numProcArrayEntries` set to zero and `plist->procArray` set to `NULL`. The return code will be either `NVML_ERROR_INSUFFICIENT_SIZE` (if there are valid processes of type `plist->mode` to report on, in which case the `plist->numProcArrayEntries` field will indicate the required number of entries in the array) or `NVML_SUCCESS` (if no processes of type `plist->mode` exist).

The `usedGpuMemory` field returned is all of the memory used by the application. The `usedGpuCcProtectedMemory` field returned is all of the protected memory used by the application.

Keep in mind that information returned by this call is dynamic and the number of elements might change in time. Allocate more space for `plist->procArray` table in case new processes are spawned.



In MIG mode, if device handle is provided, the API returns aggregate information, only if the caller has appropriate privileges. Per-instance information can be queried by using specific MIG device handles. Querying per-instance information using MIG device handles is not supported if the device is in vGPU Host virtualization mode. Protected memory usage is currently not available in MIG mode and in windows.

## `nvmlReturn_t nvmlDeviceOnSameBoard (nvmlDevice_t device1, nvmlDevice_t device2, int *onSameBoard)`

### Parameters

#### **device1**

The first GPU device

#### **device2**

The second GPU device

#### **onSameBoard**

Reference in which to return the status. Non-zero indicates that the GPUs are on the same board.

**Returns**

- ▶ NVML\_SUCCESS if onSameBoard has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if dev1 or dev2 are invalid or onSameBoard is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this check is not supported by the device
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the either GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Check if the GPU devices are on the same physical board.

For all fully supported products.

**nvmlReturn\_t nvmlDeviceGetAPIRestriction  
(nvmlDevice\_t device, nvmlRestrictedAPI\_t apiType,  
nvmlEnableState\_t \*isRestricted)**

**Parameters****device**

The identifier of the target device

**apiType**

Target API type for this operation

**isRestricted**

Reference in which to return the current restriction NVML\_FEATURE\_ENABLED indicates that the API is root-only NVML\_FEATURE\_DISABLED indicates that the API is accessible to all users

**Returns**

- ▶ NVML\_SUCCESS if isRestricted has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, apiType incorrect or isRestricted is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device or the device does not support the feature that is being queried (E.G. Enabling/disabling Auto Boosted clocks is not supported by the device)
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

## Description

Retrieves the root/admin permissions on the target API. See `nvmlRestrictedAPI_t` for the list of supported APIs. If an API is restricted only root users can call that API. See `nvmlDeviceSetAPIRestriction` to change current permissions.

For all fully supported products.

## See also:

[nvmlRestrictedAPI\\_t](#)

`nvmlReturn_t nvmlDeviceGetSamples (nvmlDevice_t device, nvmlSamplingType_t type, unsigned long long lastSeenTimeStamp, nvmlValueType_t *sampleValType, unsigned int *sampleCount, nvmlSample_t *samples)`

## Parameters

### **device**

The identifier for the target device

### **type**

Type of sampling event

### **lastSeenTimeStamp**

Return only samples with timestamp greater than lastSeenTimeStamp.

### **sampleValType**

Output parameter to represent the type of sample value as described in `nvmlSampleVal_t`

### **sampleCount**

Reference to provide the number of elements which can be queried in samples array

### **samples**

Reference in which samples are returned

## Returns

- ▶ `NVML_SUCCESS` if samples are successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, `sampleCount` is `NULL` or reference to `sampleCount` is 0 for non null samples
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_NOT_FOUND` if sample entries are not found
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

## Description

Gets recent samples for the GPU.

For Kepler or newer fully supported devices.

Based on type, this method can be used to fetch the power, utilization or clock samples maintained in the buffer by the driver.

Power, Utilization and Clock samples are returned as type "unsigned int" for the union `nvmlValue_t`.

To get the size of samples that user needs to allocate, the method is invoked with samples set to NULL. The returned samplesCount will provide the number of samples that can be queried. The user needs to allocate the buffer with size as samplesCount \* sizeof(nvmlSample\_t).

lastSeenTimeStamp represents CPU timestamp in microseconds. Set it to 0 to fetch all the samples maintained by the underlying buffer. Set lastSeenTimeStamp to one of the timeStamps retrieved from the date of the previous query to get more recent samples.

This method fetches the number of entries which can be accommodated in the provided samples array, and the reference samplesCount is updated to indicate how many samples were actually retrieved. The advantage of using this method for samples in contrast to polling via existing methods is to get higher frequency data at lower polling cost.



On MIG-enabled GPUs, querying the following sample types, NVML\_GPU\_UTILIZATION\_SAMPLES, NVML\_MEMORY\_UTILIZATION\_SAMPLES, NVML\_ENC\_UTILIZATION\_SAMPLES and NVML\_DEC\_UTILIZATION\_SAMPLES, is not currently supported.

## `nvmlReturn_t nvmlDeviceGetBAR1MemoryInfo` (`nvmlDevice_t device`, `nvmlBAR1Memory_t *bar1Memory`)

### Parameters

#### `device`

The identifier of the target device

#### `bar1Memory`

Reference in which BAR1 memory information is returned.

### Returns

- ▶ NVML\_SUCCESS if BAR1 memory is successfully retrieved

- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, `bar1Memory` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Gets Total, Available and Used size of BAR1 memory.

BAR1 is used to map the FB (device memory) so that it can be directly accessed by the CPU or by 3rd party devices (peer-to-peer on the PCIE bus).



In MIG mode, if device handle is provided, the API returns aggregate information, only if the caller has appropriate privileges. Per-instance information can be queried by using specific MIG device handles.

For Kepler or newer fully supported devices.

```
nvmlReturn_t nvmlDeviceGetViolationStatus
(nvmlDevice_t device, nvmlPerfPolicyType_t
perfPolicyType, nvmlViolationTime_t *violTime)
```

### Parameters

#### **device**

The identifier of the target device

#### **perfPolicyType**

Represents Performance policy which can trigger GPU throttling

#### **violTime**

Reference to which violation time related information is returned

### Returns

- ▶ `NVML_SUCCESS` if violation time is successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, `perfPolicyType` is invalid, or `violTime` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible

**Description**

Gets the duration of time during which the device was throttled (lower than requested clocks) due to power or thermal constraints.

The method is important to users who are trying to understand if their GPUs throttle at any point during their applications. The difference in violation times at two different reference times gives the indication of GPU throttling event.

Violation for thermal capping is not supported at this time.

For Kepler or newer fully supported devices.

## **nvmlReturn\_t nvmlDeviceGetIrqNum (nvmlDevice\_t device, unsigned int \*irqNum)**

**Parameters****device**

The identifier of the target device

**irqNum**

The interrupt number associated with the specified device

**Returns**

- ▶ NVML\_SUCCESS if irq number is successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, or irqNum is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible

**Description**

Gets the device's interrupt number

## **nvmlReturn\_t nvmlDeviceGetNumGpuCores (nvmlDevice\_t device, unsigned int \*numCores)**

**Parameters****device**

The identifier of the target device

**numCores**

The number of cores for the specified device

**Returns**

- ▶ NVML\_SUCCESS if Gpu core count is successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, or numCores is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible

**Description**

Gets the device's core count

**nvmlReturn\_t nvmlDeviceGetPowerSource**  
**(nvmlDevice\_t device, nvmlPowerSource\_t**  
**\*powerSource)**

**Parameters****device**

The identifier of the target device

**powerSource**

The power source of the device

**Returns**

- ▶ NVML\_SUCCESS if the current power source was successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, or powerSource is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible

**Description**

Gets the devices power source



## `nvmlReturn_t nvmlDeviceGetMemoryBusWidth` (`nvmlDevice_t device`, `unsigned int *busWidth`)

### Parameters

#### **device**

The identifier of the target device

#### **busWidth**

The device's memory bus width

### Returns

- ▶ `NVML_SUCCESS` if the memory bus width is successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, or busWidth is NULL
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible

### Description

Gets the device's memory bus width

## `nvmlReturn_t nvmlDeviceGetPcieLinkMaxSpeed` (`nvmlDevice_t device`, `unsigned int *maxSpeed`)

### Parameters

#### **device**

The identifier of the target device

#### **maxSpeed**

The device's PCIE Max Link speed in MBPS

### Returns

- ▶ `NVML_SUCCESS` if Pcie Max Link Speed is successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, or maxSpeed is NULL
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible

**Description**

Gets the device's PCIe Max Link speed in MBPS

**`nvmlReturn_t nvmlDeviceGetPcieSpeed (nvmlDevice_t device, unsigned int *pcieSpeed)`**

**Parameters****device**

The identifier of the target device

**pcieSpeed**

The devices's PCIe Max Link speed in Mbps

**Returns**

- ▶ NVML\_SUCCESS if pcieSpeed has been retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or pcieSpeed is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support PCIe speed getting
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Gets the device's PCIe Link speed in Mbps

**`nvmlReturn_t nvmlDeviceGetAdaptiveClockInfoStatus (nvmlDevice_t device, unsigned int *adaptiveClockStatus)`**

**Parameters****device**

The identifier of the target device

**adaptiveClockStatus**

The current adaptive clocking status, either NVML\_ADAPTIVE\_CLOCKING\_INFO\_STATUS\_DISABLED or NVML\_ADAPTIVE\_CLOCKING\_INFO\_STATUS\_ENABLED

**Returns**

- ▶ NVML\_SUCCESS if the current adaptive clocking status is successfully retrieved

- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, or `adaptiveClockStatus` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible

### Description

Gets the device's Adaptive Clock status

## `nvmlReturn_t nvmlDeviceGetBusType (nvmlDevice_t device, nvmlBusType_t *type)`

### Parameters

#### **device**

The identifier of the target device

#### **type**

The PCI Bus type

### Description

Get the type of the GPU Bus (PCIe, PCI, ...)

return

- ▶ `NVML_SUCCESS` if the bus type is successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if is invalid or is `NULL`
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

## `nvmlReturn_t nvmlDeviceGetGpuFabricInfo (nvmlDevice_t device, nvmlGpuFabricInfo_t *gpuFabricInfo)`

### Parameters

#### **device**

The identifier of the target device

#### **gpuFabricInfo**

Information about GPU fabric state

## Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't support gpu fabric

## Description

Deprecated: Will be deprecated in a future release. Use `nvmlDeviceGetGpuFabricInfoV` instead

Get fabric information associated with the device.

HOPPER\_OR\_NEWER%

On Hopper + NVSwitch systems, GPU is registered with the NVIDIA Fabric Manager. Upon successful registration, the GPU is added to the NVLink fabric to enable peer-to-peer communication. This API reports the current state of the GPU in the NVLink fabric along with other useful information.

```
nvmlReturn_t nvmlDeviceGetGpuFabricInfoV
(nvmlDevice_t device, nvmlGpuFabricInfoV_t
*gpuFabricInfo)
```

## Parameters

### device

The identifier of the target device

### gpuFabricInfo

Information about GPU fabric state

## Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't support gpu fabric

## Description

Versioned wrapper around `nvmlDeviceGetGpuFabricInfo` that accepts a versioned `nvmlGpuFabricInfo_v2_t` or later output structure.



The caller must set the `nvmlGpuFabricInfoV_t::version` field to the appropriate version prior to calling this function. For example:

```
nvmlGpuFabricInfoV_t fabricInfo =
    { .version = nvmlGpuFabricInfo_v2 };
nvmlReturn_t result
= nvmlDeviceGetGpuFabricInfoV(device, &fabricInfo);
```

HOPPER\_OR\_NEWER%

## `nvmlReturn_t nvmlSystemGetConfComputeCapabilities` (`nvmlConfComputeSystemCaps_t *capabilities`)

### Parameters

#### `capabilities`

System CC capabilities

### Returns

- ▶ `NVML_SUCCESS` if capabilities were successfully queried
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if capabilities is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device

### Description

Get Conf Computing System capabilities.

For Ampere or newer fully supported devices. Supported on Linux, Windows TCC.

## `nvmlReturn_t nvmlSystemGetConfComputeState` (`nvmlConfComputeSystemState_t *state`)

### Parameters

#### `state`

System CC State

### Returns

- ▶ `NVML_SUCCESS` if state were successfully queried
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if state is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device

### Description

Get Conf Computing System State.

For Ampere or newer fully supported devices. Supported on Linux, Windows TCC.

## `nvmlReturn_t nvmlDeviceGetConfComputeMemSizeInfo (nvmlDevice_t device, nvmlConfComputeMemSizeInfo_t *memInfo)`

### Parameters

#### **device**

Device handle

#### **memInfo**

Protected/Unprotected Memory sizes

### Returns

- ▶ `NVML_SUCCESS` if `memInfo` were successfully queried
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `memInfo` or `device` is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device

### Description

Get Conf Computing Protected and Unprotected Memory Sizes.

For Ampere or newer fully supported devices. Supported on Linux, Windows TCC.

## `nvmlReturn_t nvmlSystemGetConfComputeGpusReadyState (unsigned int *isAcceptingWork)`

### Parameters

#### **isAcceptingWork**

Returns GPU current work accepting state,  
`NVML_CC_ACCEPTING_CLIENT_REQUESTS_TRUE` or  
`NVML_CC_ACCEPTING_CLIENT_REQUESTS_FALSE`

### Description

Get Conf Computing GPUs ready state.

For Ampere or newer fully supported devices. Supported on Linux, Windows TCC.

return

- ▶ `NVML_SUCCESS` if current GPUs ready state were successfully queried

- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `isAcceptingWork` is NULL
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device

## `nvmlReturn_t` `nvmlDeviceGetConfComputeProtectedMemoryUsage` (`nvmlDevice_t` device, `nvmlMemory_t *memory`)

### Parameters

#### **device**

The identifier of the target device

#### **memory**

Reference in which to return the memory information

### Returns

- ▶ `NVML_SUCCESS` if memory has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or memory is NULL
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Get Conf Computing protected memory usage.

For Ampere or newer fully supported devices. Supported on Linux, Windows TCC.

## `nvmlReturn_t` `nvmlDeviceGetConfComputeGpuCertificate` (`nvmlDevice_t` device, `nvmlConfComputeGpuCertificate_t *gpuCert`)

### Parameters

#### **device**

The identifier of the target device

#### **gpuCert**

Reference in which to return the gpu certificate information

**Returns**

- ▶ NVML\_SUCCESS if gpu certificate info has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or memory is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Get Conf Computing Gpu certificate details.

For Ampere or newer fully supported devices. Supported on Linux, Windows TCC.

```

nvmlReturn_t
nvmlDeviceGetConfComputeGpuAttestationReport
(nvmlDevice_t device,
nvmlConfComputeGpuAttestationReport_t
*gpuAtstReport)

```

**Parameters****device**

The identifier of the target device

**gpuAtstReport**

Reference in which to return the gpu attestation report

**Returns**

- ▶ NVML\_SUCCESS if gpu attestation report has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or memory is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Get Conf Computing Gpu attestation report.

For Ampere or newer fully supported devices. Supported on Linux, Windows TCC.

```

nvmlReturn_t
nvmlSystemGetConfComputeKeyRotationThresholdInfo

```



## (nvmlConfComputeGetKeyRotationThresholdInfo\_t \*pKeyRotationThrInfo)

### Parameters

#### pKeyRotationThrInfo

Reference in which to return the key rotation threshold data

### Returns

- ▶ NVML\_SUCCESS if gpu key rotation threshold info has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or memory is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Get Conf Computing key rotation threshold detail.

HOPPER\_OR\_NEWER% Supported on Linux, Windows TCC.

## nvmlReturn\_t nvmlSystemGetConfComputeSettings (nvmlSystemConfComputeSettings\_t \*settings)

### Parameters

#### settings

System CC settings

### Returns

- ▶ NVML\_SUCCESS if the query is success
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or counters is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_ARGUMENT\_VERSION\_MISMATCH if the provided version is invalid/unsupported
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Get Conf Computing System Settings.

HOPPER\_OR\_NEWER% Supported on Linux, Windows TCC.

## **nvmlReturn\_t nvmlDeviceGetGspFirmwareVersion (nvmlDevice\_t device, char \*version)**

**Parameters****device**

Device handle

**version**

The retrieved GSP firmware version

**Returns**

- ▶ NVML\_SUCCESS if GSP firmware version is successfully retrieved
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or GSP version pointer is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if GSP firmware is not enabled for GPU
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve GSP firmware version.

The caller passes in buffer via version and corresponding GSP firmware numbered version is returned with the same parameter in string format.

## **nvmlReturn\_t nvmlDeviceGetGspFirmwareMode (nvmlDevice\_t device, unsigned int \*isEnabled, unsigned int \*defaultMode)**

**Parameters****device**

Device handle

**isEnabled**

Pointer to specify if GSP firmware is enabled

**defaultMode**

Pointer to specify if GSP firmware is supported by default on device

**Returns**

- ▶ NVML\_SUCCESS if GSP firmware mode is successfully retrieved
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or any of isEnabled or defaultMode is NULL
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve GSP firmware mode.

The caller passes in integer pointers. GSP firmware enablement and default mode information is returned with corresponding parameters. The return value in isEnabled and defaultMode should be treated as boolean.

**nvmlReturn\_t nvmlDeviceGetRetiredPages**  
**(nvmlDevice\_t device, nvmlPageRetirementCause\_t**  
**cause, unsigned int \*pageCount, unsigned long long**  
**\*addresses)**

**Parameters****device**

The identifier of the target device

**cause**

Filter page addresses by cause of retirement

**pageCount**

Reference in which to provide the addresses buffer size, and to return the number of retired pages that match cause Set to 0 to query the size without allocating an addresses buffer

**addresses**

Buffer to write the page addresses into

**Returns**

- ▶ NVML\_SUCCESS if pageCount was populated and addresses was filled
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if pageCount indicates the buffer is not large enough to store all the matching page addresses. pageCount is set to the needed size.
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, pageCount is NULL, cause is invalid, or addresses is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature

- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Returns the list of retired pages by source, including pages that are pending retirement. The address information provided from this API is the hardware address of the page that was retired. Note that this does not match the virtual address used in CUDA, but will match the address information in `XID 63`.

For Kepler or newer fully supported devices.

`nvmlReturn_t nvmlDeviceGetRetiredPages_v2`  
 (`nvmlDevice_t device`, `nvmlPageRetirementCause_t cause`, `unsigned int *pageCount`, `unsigned long long *addresses`, `unsigned long long *timestamps`)

### Parameters

#### `device`

The identifier of the target device

#### `cause`

Filter page addresses by cause of retirement

#### `pageCount`

Reference in which to provide the addresses buffer size, and to return the number of retired pages that match `cause`. Set to 0 to query the size without allocating an addresses buffer.

#### `addresses`

Buffer to write the page addresses into

#### `timestamps`

Buffer to write the timestamps of page retirement, additional for `_v2`

### Returns

- ▶ `NVML_SUCCESS` if `pageCount` was populated and `addresses` was filled
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `pageCount` indicates the buffer is not large enough to store all the matching page addresses. `pageCount` is set to the needed size.
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid, `pageCount` is `NULL`, `cause` is invalid, or `addresses` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device doesn't support this feature

- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Returns the list of retired pages by source, including pages that are pending retirement. The address information provided from this API is the hardware address of the page that was retired. Note that this does not match the virtual address used in CUDA, but will match the address information in `XID 63`.



`nvmlDeviceGetRetiredPages_v2` adds an additional `timestamps` parameter to return the time of each page's retirement.

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetRetiredPagesPendingStatus` (`nvmlDevice_t device, nvmlEnableState_t *isPending`)

### Parameters

#### **device**

The identifier of the target device

#### **isPending**

Reference in which to return the pending status

### Returns

- ▶ `NVML_SUCCESS` if `isPending` was populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid or `isPending` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device doesn't support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Check if any pages are pending retirement and need a reboot to fully retire.

For Kepler or newer fully supported devices.

**nvmlReturn\_t nvmlDeviceGetRemappedRows**  
 (nvmlDevice\_t device, unsigned int \*corrRows, unsigned int \*uncRows, unsigned int \*isPending, unsigned int \*failureOccurred)

### Parameters

#### **device**

The identifier of the target device

#### **corrRows**

Reference for number of rows remapped due to correctable errors

#### **uncRows**

Reference for number of rows remapped due to uncorrectable errors

#### **isPending**

Reference for whether or not remappings are pending

#### **failureOccurred**

Reference that is set when a remapping has failed in the past

### Returns

- ▶ NVML\_SUCCESS Upon success
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT If corrRows, uncRows, isPending or failureOccurred is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED If MIG is enabled or if the device doesn't support this feature
- ▶ NVML\_ERROR\_UNKNOWN Unexpected error

### Description

Get number of remapped rows. The number of rows reported will be based on the cause of the remapping. isPending indicates whether or not there are pending remappings. A reset will be required to actually remap the row. failureOccurred will be set if a row remapping ever failed in the past. A pending remapping won't affect future work on the GPU since error-containment and dynamic page blacklisting will take care of that.



On MIG-enabled GPUs with active instances, querying the number of remapped rows is not supported

For Ampere or newer fully supported devices.

**nvmlReturn\_t nvmlDeviceGetRowRemapperHistogram**  
**(nvmlDevice\_t device,**  
**nvmlRowRemapperHistogramValues\_t \*values)**

#### Parameters

##### **device**

Device handle

##### **values**

Histogram values

#### Returns

- ▶ NVML\_SUCCESS On success
- ▶ NVML\_ERROR\_UNKNOWN On any unexpected error

#### Description

Get the row remapper histogram. Returns the remap availability for each bank on the GPU.

**nvmlReturn\_t nvmlDeviceGetArchitecture**  
**(nvmlDevice\_t device, nvmlDeviceArchitecture\_t \*arch)**

#### Parameters

##### **device**

The identifier of the target device

##### **arch**

Reference where architecture is returned, if call successful. Set to NVML\_DEVICE\_ARCH\_\* upon success

#### Returns

- ▶ NVML\_SUCCESS Upon success
- ▶ NVML\_ERROR\_UNINITIALIZED If library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT If device or arch (output reference) are invalid

#### Description

Get architecture for device

## `nvmlReturn_t nvmlDeviceGetClkMonStatus` (`nvmlDevice_t device`, `nvmlClkMonStatus_t *status`)

### Parameters

#### **device**

The identifier of the target device

#### **status**

Reference in which to return the clkmon fault status

### Returns

- ▶ `NVML_SUCCESS` if status has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or status is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the frequency monitor fault status for the device.

For Ampere or newer fully supported devices. Requires root user.

See `nvmlClkMonStatus_t` for details on decoding the status output.

### See also:

`nvmlDeviceGetClkMonStatus()`

## `nvmlReturn_t nvmlDeviceGetProcessUtilization` (`nvmlDevice_t device`, `nvmlProcessUtilizationSample_t *utilization`, `unsigned int *processSamplesCount`, `unsigned long long lastSeenTimeStamp`)

### Parameters

#### **device**

The identifier of the target device



**utilization**

Pointer to caller-supplied buffer in which guest process utilization samples are returned

**processSamplesCount**

Pointer to caller-supplied array size, and returns number of processes running

**lastSeenTimeStamp**

Return only samples with timestamp greater than lastSeenTimeStamp.

**Returns**

- ▶ NVML\_SUCCESS if utilization has been populated
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, utilization is NULL, or samplingPeriodUs is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_NOT\_FOUND if sample entries are not found
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the current utilization and process ID

For Maxwell or newer fully supported devices.

Reads recent utilization of GPU SM (3D/Compute), framebuffer, video encoder, and video decoder for processes running. Utilization values are returned as an array of utilization sample structures in the caller-supplied buffer pointed at by utilization. One utilization sample structure is returned per process running, that had some non-zero utilization during the last sample period. It includes the CPU timestamp at which the samples were recorded. Individual utilization values are returned as "unsigned int" values. If no valid sample entries are found since the lastSeenTimeStamp, NVML\_ERROR\_NOT\_FOUND is returned.

To read utilization values, first determine the size of buffer required to hold the samples by invoking the function with utilization set to NULL. The caller should allocate a buffer of size processSamplesCount \* sizeof(nvmlProcessUtilizationSample\_t). Invoke the function again with the allocated buffer passed in utilization, and processSamplesCount set to the number of entries the buffer is sized for.

On successful return, the function updates processSamplesCount with the number of process utilization sample structures that were actually written. This may differ from a previously read value as instances are created or destroyed.

lastSeenTimeStamp represents the CPU timestamp in microseconds at which utilization samples were last read. Set it to 0 to read utilization based on all the samples maintained

by the driver's internal sample buffer. Set `lastSeenTimeStamp` to a `timeStamp` retrieved from a previous query to read utilization since the previous query.



On MIG-enabled GPUs, querying process utilization is not currently supported.

## `nvmlReturn_t nvmlDeviceGetProcessesUtilizationInfo` (`nvmlDevice_t device`, `nvmlProcessesUtilizationInfo_t` `*procesesUtilInfo`)

### Parameters

#### `device`

The identifier of the target device

#### `procesesUtilInfo`

Pointer to the caller-provided structure of `nvmlProcessesUtilizationInfo_t`.

### Returns

- ▶ `NVML_SUCCESS` if `procesesUtilInfo->procUtilArray` has been populated
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid, or `procesesUtilInfo` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_NOT_FOUND` if sample entries are not found
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_VERSION_MISMATCH` if the version of `procesesUtilInfo` is invalid
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `procesesUtilInfo->procUtilArray` is `NULL`, or the buffer size of `procesesUtilInfo->procUtilArray` is too small. The caller should check the minimum array size from the returned `procesesUtilInfo->processSamplesCount`, and call the function again with a buffer no smaller than `procesesUtilInfo->processSamplesCount * sizeof(nvmlProcessUtilizationInfo_t)`
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the recent utilization and process ID for all running processes

For Maxwell or newer fully supported devices.

Reads recent utilization of GPU SM (3D/Compute), framebuffer, video encoder, and video decoder, jpeg decoder, OFA (Optical Flow Accelerator) for all running processes. Utilization values are returned as an array of utilization sample structures in the caller-

supplied buffer pointed at by `procesesUtilInfo->procUtilArray`. One utilization sample structure is returned per process running, that had some non-zero utilization during the last sample period. It includes the CPU timestamp at which the samples were recorded. Individual utilization values are returned as "unsigned int" values.

The caller should allocate a buffer of size `processSamplesCount * sizeof(nvmlProcessUtilizationInfo_t)`. If the buffer is too small, the API will return `NVML_ERROR_INSUFFICIENT_SIZE`, with the recommended minimal buffer size at `procesesUtilInfo->processSamplesCount`. The caller should invoke the function again with the allocated buffer passed in `procesesUtilInfo->procUtilArray`, and `procesesUtilInfo->processSamplesCount` set to the number no less than the recommended value by the previous API return.

On successful return, the function updates `procesesUtilInfo->processSamplesCount` with the number of process utilization info structures that were actually written. This may differ from a previously read value as instances are created or destroyed.

`procesesUtilInfo->lastSeenTimeStamp` represents the CPU timestamp in microseconds at which utilization samples were last read. Set it to 0 to read utilization based on all the samples maintained by the driver's internal sample buffer. Set `procesesUtilInfo->lastSeenTimeStamp` to a `timeStamp` retrieved from a previous query to read utilization since the previous query.

`procesesUtilInfo->version` is the version number of the structure `nvmlProcessesUtilizationInfo_t`, the caller should set the correct version number to retrieve the specific version of processes utilization information.



On MIG-enabled GPUs, querying process utilization is not currently supported.

## 4.15.1. CPU and Memory Affinity

### Device Queries

This chapter describes NVML operations that are associated with CPU and memory affinity.

```
nvmlReturn_t nvmlDeviceGetMemoryAffinity (nvmlDevice_t
device, unsigned int nodeSetSize, unsignedlong *nodeSet,
nvmlAffinityScope_t scope)
```

### Parameters

#### **device**

The identifier of the target device

**nodeSetSize**

The size of the nodeSet array that is safe to access

**nodeSet**

Array reference in which to return a bitmask of NODEs, 64 NODEs per unsigned long on 64-bit machines, 32 on 32-bit machines

**scope**

Scope that change the default behavior

**Returns**

- ▶ NVML\_SUCCESS if NUMA node Affinity has been filled
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, nodeSetSize == 0, nodeSet is NULL or scope is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves an array of unsigned ints (sized to nodeSetSize) of bitmasks with the ideal memory affinity within node or socket for the device. For example, if NUMA node 0, 1 are ideal within the socket for the device and nodeSetSize == 1, result[0] = 0x3



If requested scope is not applicable to the target topology, the API will fall back to reporting the memory affinity for the immediate non-I/O ancestor of the device.

For Kepler or newer fully supported devices. Supported on Linux only.

**nvmlReturn\_t nvmlDeviceGetCpuAffinityWithinScope (nvmlDevice\_t device, unsigned int cpuSetSize, unsignedlong \*cpuSet, nvmlAffinityScope\_t scope)**

**Parameters****device**

The identifier of the target device

**cpuSetSize**

The size of the cpuSet array that is safe to access

**cpuSet**

Array reference in which to return a bitmask of CPUs, 64 CPUs per unsigned long on 64-bit machines, 32 on 32-bit machines

**scope**

Scope that change the default behavior

**Returns**

- ▶ NVML\_SUCCESS if cpuAffinity has been filled
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, cpuSetSize == 0, cpuSet is NULL or scope is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves an array of unsigned ints (sized to cpuSetSize) of bitmasks with the ideal CPU affinity within node or socket for the device. For example, if processors 0, 1, 32, and 33 are ideal for the device and cpuSetSize == 2, result[0] = 0x3, result[1] = 0x3



If requested scope is not applicable to the target topology, the API will fall back to reporting the CPU affinity for the immediate non-I/O ancestor of the device.

For Kepler or newer fully supported devices. Supported on Linux only.

**`nvmlReturn_t nvmlDeviceGetCpuAffinity (nvmlDevice_t device, unsigned int cpuSetSize, unsignedlong *cpuSet)`**

**Parameters****device**

The identifier of the target device

**cpuSetSize**

The size of the cpuSet array that is safe to access

**cpuSet**

Array reference in which to return a bitmask of CPUs, 64 CPUs per unsigned long on 64-bit machines, 32 on 32-bit machines

**Returns**

- ▶ NVML\_SUCCESS if cpuAffinity has been filled
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, cpuSetSize == 0, or cpuSet is NULL

- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves an array of unsigned ints (sized to `cpuSetSize`) of bitmasks with the ideal CPU affinity for the device. For example, if processors 0, 1, 32, and 33 are ideal for the device and `cpuSetSize == 2`, `result[0] = 0x3`, `result[1] = 0x3`. This is equivalent to calling [`nvmlDeviceGetCpuAffinityWithinScope`](#) with `NVML_AFFINITY_SCOPE_NODE`.

For Kepler or newer fully supported devices. Supported on Linux only.

## `nvmlReturn_t nvmlDeviceSetCpuAffinity (nvmlDevice_t device)`

### Parameters

#### `device`

The identifier of the target device

### Returns

- ▶ `NVML_SUCCESS` if the calling process has been successfully bound
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Sets the ideal affinity for the calling thread and device using the guidelines given in [`nvmlDeviceGetCpuAffinity\(\)`](#). Note, this is a change as of version 8.0. Older versions set the affinity for a calling process and all children. Currently supports up to 1024 processors.

For Kepler or newer fully supported devices. Supported on Linux only.

## `nvmlReturn_t nvmlDeviceClearCpuAffinity (nvmlDevice_t device)`

### Parameters

#### **device**

The identifier of the target device

### Returns

- ▶ `NVML_SUCCESS` if the calling process has been successfully unbound
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Clear all affinity bindings for the calling thread. Note, this is a change as of version 8.0 as older versions cleared the affinity for a calling process and all children.

For Kepler or newer fully supported devices. Supported on Linux only.

## `nvmlReturn_t nvmlDeviceGetNumaNodeid (nvmlDevice_t device, unsigned int *node)`

### Parameters

#### **device**

The device handle

#### **node**

NUMA node ID of the device

### Returns

- ▶ `NVML_SUCCESS` if the NUMA node is retrieved successfully
- ▶ `NVML_ERROR_NOT_SUPPORTED` if request is not supported on the current platform
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device node is invalid

### Description

Get the NUMA node of the given GPU device. This only applies to platforms where the GPUs are NUMA nodes.

```
#define NVML_AFFINITY_SCOPE_NODE 0
```

Scope of NUMA node for affinity queries.

```
#define NVML_AFFINITY_SCOPE_SOCKET 1
```

Scope of processor socket for affinity queries.

## 4.16. Unit Commands

This chapter describes NVML operations that change the state of the unit. For S-class products. Each of these requires root/admin access. Non-admin users will see an NVML\_ERROR\_NO\_PERMISSION error code when invoking any of these methods.

```
nvmlReturn_t nvmlUnitSetLedState (nvmlUnit_t unit,
nvmlLedColor_t color)
```

### Parameters

#### unit

The identifier of the target unit

#### color

The target LED color

### Returns

- ▶ NVML\_SUCCESS if the LED color has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if unit or color is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this is not an S-class product
- ▶ NVML\_ERROR\_NO\_PERMISSION if the user doesn't have permission to perform this operation
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Set the LED state for the unit. The LED can be either green (0) or amber (1).

For S-class products. Requires root/admin permissions.

This operation takes effect immediately.

**Current S-Class products don't provide unique LEDs for each unit. As such, both front and back LEDs will be toggled in unison regardless of which unit is specified with this command.**



See `nvmlLedColor_t` for available colors.

**See also:**

`nvmlUnitGetLedState()`

## 4.17. Device Commands

This chapter describes NVML operations that change the state of the device. Each of these requires root/admin access. Non-admin users will see an `NVML_ERROR_NO_PERMISSION` error code when invoking any of these methods.

### `nvmlReturn_t nvmlDeviceSetPersistenceMode` (`nvmlDevice_t device`, `nvmlEnableState_t mode`)

#### Parameters

##### **device**

The identifier of the target device

##### **mode**

The target persistence mode

#### Returns

- ▶ `NVML_SUCCESS` if the persistence mode was set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or mode is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

#### Description

Set the persistence mode for the device.

For all products. For Linux only. Requires root/admin permissions.

The persistence mode determines whether the GPU driver software is torn down after the last client exits.

This operation takes effect immediately. It is not persistent across reboots. After each reboot the persistence mode is reset to "Disabled".

See [nvmlEnableState\\_t](#) for available modes.

After calling this API with mode set to `NVML_FEATURE_DISABLED` on a device that has its own NUMA memory, the given device handle will no longer be valid, and to continue to interact with this device, a new handle should be obtained from one of the `nvmlDeviceGetHandleBy*`() APIs. This limitation is currently only applicable to devices that have a coherent NVLink connection to system memory.

**See also:**

[nvmlDeviceGetPersistenceMode\(\)](#)

## `nvmlReturn_t nvmlDeviceSetComputeMode` (`nvmlDevice_t device`, `nvmlComputeMode_t mode`)

### Parameters

#### **device**

The identifier of the target device

#### **mode**

The target compute mode

### Returns

- ▶ `NVML_SUCCESS` if the compute mode was set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or mode is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Set the compute mode for the device.

For all products. Requires root/admin permissions.

The compute mode determines whether a GPU can be used for compute operations and whether it can be shared across contexts.

This operation takes effect immediately. Under Linux it is not persistent across reboots and always resets to "Default". Under windows it is persistent.

Under windows compute mode may only be set to DEFAULT when running in WDDM



On MIG-enabled GPUs, compute mode would be set to DEFAULT and changing it is not supported.

See `nvmlComputeMode_t` for details on available compute modes.

#### See also:

`nvmlDeviceGetComputeMode()`

## `nvmlReturn_t nvmlDeviceSetEccMode (nvmlDevice_t device, nvmlEnableState_t ecc)`

### Parameters

#### `device`

The identifier of the target device

#### `ecc`

The target ECC mode

### Returns

- ▶ `NVML_SUCCESS` if the ECC mode was set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or ecc is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Set the ECC mode for the device.

For Kepler or newer fully supported devices. Only applicable to devices with ECC. Requires `NVML_INFOROM_ECC` version 1.0 or higher. Requires root/admin permissions.

The ECC mode determines whether the GPU enables its ECC support.

This operation takes effect after the next reboot.

See [nvmlEnableState\\_t](#) for details on available modes.

**See also:**

[nvmlDeviceGetEccMode\(\)](#)

## **`nvmlReturn_t nvmlDeviceClearEccErrorCounts` (`nvmlDevice_t device`, `nvmlEccCounterType_t counterType`)**

**Parameters**

**device**

The identifier of the target device

**counterType**

Flag that indicates which type of errors should be cleared.

**Returns**

- ▶ NVML\_SUCCESS if the error counts were cleared
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or counterType is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_NO\_PERMISSION if the user doesn't have permission to perform this operation
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Clear the ECC error and other memory error counts for the device.

For Kepler or newer fully supported devices. Only applicable to devices with ECC. Requires NVML\_INFOROM\_ECC version 2.0 or higher to clear aggregate location-based ECC counts. Requires NVML\_INFOROM\_ECC version 1.0 or higher to clear all other ECC counts. Requires root/admin permissions. Requires ECC Mode to be enabled.

Sets all of the specified ECC counters to 0, including both detailed and total counts.

This operation takes effect immediately.

See [nvmlMemoryErrorType\\_t](#) for details on available counter types.

**See also:**

- ▶ [nvmlDeviceGetDetailedEccErrors\(\)](#)
- ▶ [nvmlDeviceGetTotalEccErrors\(\)](#)

## **nvmlReturn\_t nvmlDeviceSetDriverModel (nvmlDevice\_t device, nvmlDriverModel\_t driverModel, unsigned int flags)**

**Parameters****device**

The identifier of the target device

**driverModel**

The target driver model

**flags**

Flags that change the default behavior

**Returns**

- ▶ NVML\_SUCCESS if the driver model has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or driverModel is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the platform is not windows or the device does not support this feature
- ▶ NVML\_ERROR\_NO\_PERMISSION if the user doesn't have permission to perform this operation
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Set the driver model for the device.

For Fermi or newer fully supported devices. For windows only. Requires root/admin permissions.

On Windows platforms the device driver can run in either WDDM or WDM (TCC) mode. If a display is attached to the device it must run in WDDM mode.

It is possible to force the change to WDM (TCC) while the display is still attached with a force flag (nvmlFlagForce). This should only be done if the host is subsequently powered down and the display is detached from the device before the next reboot.

This operation takes effect after the next reboot.

Windows driver model may only be set to WDDM when running in DEFAULT compute mode.

Change driver model to WDDM is not supported when GPU doesn't support graphics acceleration or will not support it after reboot. See [nvmlDeviceSetGpuOperationMode](#).

See [nvmlDriverModel\\_t](#) for details on available driver models. See [nvmlFlagDefault](#) and [nvmlFlagForce](#)

#### See also:

[nvmlDeviceGetDriverModel\(\)](#)

## **`nvmlReturn_t nvmlDeviceSetGpuLockedClocks`** **(`nvmlDevice_t device`, `unsigned int minGpuClockMHz`, `unsigned int maxGpuClockMHz`)**

### Parameters

#### **device**

The identifier of the target device

#### **minGpuClockMHz**

Requested minimum gpu clock in MHz

#### **maxGpuClockMHz**

Requested maximum gpu clock in MHz

### Returns

- ▶ `NVML_SUCCESS` if new settings were successfully set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or `minGpuClockMHz` and `maxGpuClockMHz` is not a valid clock combination
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device doesn't support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Set clocks that device will lock to.

Sets the clocks that the device will be running at to the value in the range of `minGpuClockMHz` to `maxGpuClockMHz`. Setting this will supersede application clock values and take effect regardless if a cuda app is running. See [/ref nvidiaDeviceSetApplicationsClocks](#)

Can be used as a setting to request constant performance.

This can be called with a pair of integer clock frequencies in MHz, or a pair of [/ref nvidiaClockLimitId\\_t](#) values. See the table below for valid combinations of these values.

<code>minGpuClock</code>	<code>maxGpuClock</code>	Effect
tdp	tdp	Lock clock to TDP unlimited
tdp	unlimited	Upper bound is TDP but clock may drift below this
unlimited	tdp	Lower bound is TDP but clock may boost above this
unlimited	unlimited	Unlocked (= <code>nvidiaDeviceResetGpuLockedClocks</code> )

If one arg takes one of these values, the other must be one of these values as well. Mixed numeric and symbolic calls return `NVML_ERROR_INVALID_ARGUMENT`.

Requires root/admin permissions.

After system reboot or driver reload applications clocks go back to their default value. See [nvidiaDeviceResetGpuLockedClocks](#).

For Volta or newer fully supported devices.

## `nvidiaReturn_t nvidiaDeviceResetGpuLockedClocks(nvidiaDevice_t device)`

### Parameters

#### `device`

The identifier of the target device

### Returns

- ▶ `NVML_SUCCESS` if new settings were successfully set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Resets the gpu clock to the default value

This is the gpu clock that will be used after system reboot or driver reload. Default values are idle clocks, but the current values can be changed using [nvmlDeviceSetApplicationsClocks](#).

**See also:**

[nvmlDeviceSetGpuLockedClocks](#)

For Volta or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceSetMemoryLockedClocks` (`nvmlDevice_t device`, `unsigned int minMemClockMHz`, `unsigned int maxMemClockMHz`)

### Parameters

**device**

The identifier of the target device

**minMemClockMHz**

Requested minimum memory clock in MHz

**maxMemClockMHz**

Requested maximum memory clock in MHz

### Returns

- ▶ `NVML_SUCCESS` if new settings were successfully set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or `minGpuClockMHz` and `maxGpuClockMHz` is not a valid clock combination
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device doesn't support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Set memory clocks that device will lock to.

Sets the device's memory clocks to the value in the range of `minMemClockMHz` to `maxMemClockMHz`. Setting this will supersede application clock values and take effect regardless of whether a cuda app is running. See [/ref nvmlDeviceSetApplicationsClocks](#)

Can be used as a setting to request constant performance.



Requires root/admin permissions.

After system reboot or driver reload applications clocks go back to their default value. See [nvmIDeviceResetMemoryLockedClocks](#).

For Ampere or newer fully supported devices.

## `nvmIReturn_t nvmIDeviceResetMemoryLockedClocks` (`nvmIDevice_t device`)

### Parameters

#### `device`

The identifier of the target device

### Returns

- ▶ `NVML_SUCCESS` if new settings were successfully set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Resets the memory clock to the default value

This is the memory clock that will be used after system reboot or driver reload. Default values are idle clocks, but the current values can be changed using [nvmIDeviceSetApplicationsClocks](#).

#### See also:

[nvmIDeviceSetMemoryLockedClocks](#)

For Ampere or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceSetApplicationsClocks` (`nvmlDevice_t device`, `unsigned int memClockMHz`, `unsigned int graphicsClockMHz`)

### Parameters

#### `device`

The identifier of the target device

#### `memClockMHz`

Requested memory clock in MHz

#### `graphicsClockMHz`

Requested graphics clock in MHz

### Returns

- ▶ `NVML_SUCCESS` if new settings were successfully set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or `memClockMHz` and `graphicsClockMHz` is not a valid clock combination
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device doesn't support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Set clocks that applications will lock to.

Sets the clocks that compute and graphics applications will be running at. e.g. CUDA driver requests these clocks during context creation which means this property defines clocks at which CUDA applications will be running unless some overspec event occurs (e.g. over power, over thermal or external HW brake).

Can be used as a setting to request constant performance.

On Pascal and newer hardware, this will automatically disable automatic boosting of clocks.

On K80 and newer Kepler and Maxwell GPUs, users desiring fixed performance should also call `nvmlDeviceSetAutoBoostedClocksEnabled` to prevent clocks from automatically boosting above the clock value being set.

For Kepler or newer non-GeForce fully supported devices and Maxwell or newer GeForce devices. Requires root/admin permissions.

See [nvmlDeviceGetSupportedMemoryClocks](#) and [nvmlDeviceGetSupportedGraphicsClocks](#) for details on how to list available clocks combinations.

After system reboot or driver reload applications clocks go back to their default value. See [nvmlDeviceResetApplicationsClocks](#).

## **`nvmlReturn_t nvmlDeviceResetApplicationsClocks`** **(`nvmlDevice_t` device)**

### **Parameters**

#### **device**

The identifier of the target device

### **Returns**

- ▶ `NVML_SUCCESS` if new settings were successfully set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### **Description**

Resets the application clock to the default value

This is the applications clock that will be used after system reboot or driver reload. Default value is constant, but the current value can be changed using [nvmlDeviceSetApplicationsClocks](#).

On Pascal and newer hardware, if clocks were previously locked with [nvmlDeviceSetApplicationsClocks](#), this call will unlock clocks. This returns clocks their default behavior of automatically boosting above base clocks as thermal limits allow.

### **See also:**

[nvmlDeviceGetApplicationsClock](#)

[nvmlDeviceSetApplicationsClocks](#)

For Fermi or newer non-GeForce fully supported devices and Maxwell or newer GeForce devices.

## `nvmlReturn_t nvmlDeviceSetAutoBoostedClocksEnabled(nvmlDevice_t device, nvmlEnableState_t enabled)`

### Parameters

#### **device**

The identifier of the target device

#### **enabled**

What state to try to set Auto Boosted clocks of the target device to

### Returns

- ▶ `NVML_SUCCESS` If the Auto Boosted clocks were successfully set to the state specified by `enabled`
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support Auto Boosted clocks
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Try to set the current state of Auto Boosted clocks on a device.

For Kepler or newer fully supported devices.

Auto Boosted clocks are enabled by default on some hardware, allowing the GPU to run at higher clock rates to maximize performance as thermal limits allow. Auto Boosted clocks should be disabled if fixed clock rates are desired.

Non-root users may use this API by default but can be restricted by root from using this API by calling `nvmlDeviceSetAPIRestriction` with `apiType=NVML_RESTRICTED_API_SET_AUTO_BOOSTED_CLOCKS`. Note: Persistence Mode is required to modify current Auto Boost settings, therefore, it must be enabled.

On Pascal and newer hardware, Auto Boosted clocks are controlled through application clocks. Use `nvmlDeviceSetApplicationsClocks` and `nvmlDeviceResetApplicationsClocks` to control Auto Boost behavior.

## `nvmlReturn_t` `nvmlDeviceSetDefaultAutoBoostedClocksEnabled` (`nvmlDevice_t` device, `nvmlEnableState_t` enabled, unsigned int flags)

### Parameters

#### **device**

The identifier of the target device

#### **enabled**

What state to try to set default Auto Boosted clocks of the target device to

#### **flags**

Flags that change the default behavior. Currently Unused.

### Returns

- ▶ `NVML_SUCCESS` If the Auto Boosted clock's default state was successfully set to the state specified by `enabled`
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_NO_PERMISSION` If the calling user does not have permission to change Auto Boosted clock's default state.
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support Auto Boosted clocks
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Try to set the default state of Auto Boosted clocks on a device. This is the default state that Auto Boosted clocks will return to when no compute running processes (e.g. CUDA application which have an active context) are running

For Kepler or newer non-GeForce fully supported devices and Maxwell or newer GeForce devices. Requires root/admin permissions.

Auto Boosted clocks are enabled by default on some hardware, allowing the GPU to run at higher clock rates to maximize performance as thermal limits allow. Auto Boosted clocks should be disabled if fixed clock rates are desired.

On Pascal and newer hardware, Auto Boosted clocks are controlled through application clocks. Use `nvmlDeviceSetApplicationsClocks` and `nvmlDeviceResetApplicationsClocks` to control Auto Boost behavior.

## `nvmlReturn_t nvmlDeviceSetDefaultFanSpeed_v2` (`nvmlDevice_t device`, `unsigned int fan`)

### Parameters

#### **device**

The identifier of the target device

#### **fan**

The index of the fan, starting at zero

### Description

Sets the speed of the fan control policy to default.

For all cuda-capable discrete products with fans

return `NVML_SUCCESS` if speed has been adjusted

`NVML_ERROR_UNINITIALIZED` if the library has not been successfully

initialized `NVML_ERROR_INVALID_ARGUMENT` if device is invalid

`NVML_ERROR_NOT_SUPPORTED` if the device does not support this (doesn't have fans) `NVML_ERROR_UNKNOWN` on any unexpected error

## `nvmlReturn_t nvmlDeviceSetFanControlPolicy` (`nvmlDevice_t device`, `unsigned int fan`, `nvmlFanControlPolicy_t policy`)

### Description

Sets current fan control policy.

For Maxwell or newer fully supported devices.

Requires privileged user.

For all cuda-capable discrete products with fans

`device` The identifier of the target device `policy` The fan control policy to set

return `NVML_SUCCESS` if policy has been set `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized `NVML_ERROR_INVALID_ARGUMENT`

if device is invalid or policy is null or the fan given doesn't reference a fan that exists. `NVML_ERROR_NOT_SUPPORTED` if the device is older than Maxwell

`NVML_ERROR_UNKNOWN` on any unexpected error

## `nvmlReturn_t nvmlDeviceSetTemperatureThreshold` (`nvmlDevice_t device`, `nvmlTemperatureThresholds_t thresholdType`, `int *temp`)

### Parameters

#### **device**

The identifier of the target device

#### **thresholdType**

The type of threshold value to be set

#### **temp**

Reference which hold the value to be set

### Returns

- ▶ `NVML_SUCCESS` if temp has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, `thresholdType` is invalid or temp is NULL
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not have a temperature sensor or is unsupported
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Sets the temperature threshold for the GPU with the specified threshold type in degrees C.

For Maxwell or newer fully supported devices.

See `nvmlTemperatureThresholds_t` for details on available temperature thresholds.

## `nvmlReturn_t nvmlDeviceSetPowerManagementLimit` (`nvmlDevice_t device`, `unsigned int limit`)

### Parameters

#### **device**

The identifier of the target device

#### **limit**

Power management limit in milliwatts to set

**Returns**

- ▶ NVML\_SUCCESS if limit has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or defaultLimit is out of range
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device does not support this feature
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Set new power limit of this device.

For Kepler or newer fully supported devices. Requires root/admin permissions.

See [nvmlDeviceGetPowerManagementLimitConstraints](#) to check the allowed ranges of values.



Limit is not persistent across reboots or driver unloads. Enable persistent mode to prevent driver from unloading when no application is using the device.

**See also:**

[nvmlDeviceGetPowerManagementLimitConstraints](#)

[nvmlDeviceGetPowerManagementDefaultLimit](#)

## **`nvmlReturn_t nvmlDeviceSetGpuOperationMode` (`nvmlDevice_t device`, `nvmlGpuOperationMode_t mode`)**

**Parameters****device**

The identifier of the target device

**mode**

Target GOM

**Returns**

- ▶ NVML\_SUCCESS if mode has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or mode incorrect



- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support GOM or specific mode
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Sets new GOM. See `nvmlGpuOperationMode_t` for details.

For GK110 M-class and X-class Tesla products from the Kepler family. Modes `NVML_GOM_LOW_DP` and `NVML_GOM_ALL_ON` are supported on fully supported GeForce products. Not supported on Quadro and Tesla C-class products. Requires root/admin permissions.

Changing GOMs requires a reboot. The reboot requirement might be removed in the future.

Compute only GOMs don't support graphics acceleration. Under windows switching to these GOMs when pending driver model is WDDM is not supported. See `nvmlDeviceSetDriverModel`.

### See also:

`nvmlGpuOperationMode_t`

`nvmlDeviceGetGpuOperationMode`

`nvmlReturn_t nvmlDeviceSetAPIRestriction`  
(`nvmlDevice_t device`, `nvmlRestrictedAPI_t apiType`,  
`nvmlEnableState_t isRestricted`)

### Parameters

#### **device**

The identifier of the target device

#### **apiType**

Target API type for this operation

#### **isRestricted**

The target restriction

### Returns

- ▶ `NVML_SUCCESS` if `isRestricted` has been set

- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or `apiType` incorrect
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support changing API restrictions or the device does not support the feature that api restrictions are being set for (E.G. Enabling/disabling auto boosted clocks is not supported by the device)
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Changes the root/admin restrictions on certain APIs. See `nvmlRestrictedAPI_t` for the list of supported APIs. This method can be used by a root/admin user to give non-root/admin access to certain otherwise-restricted APIs. The new setting lasts for the lifetime of the NVIDIA driver; it is not persistent. See `nvmlDeviceGetAPIRestriction` to query the current restriction settings.

For Kepler or newer fully supported devices. Requires root/admin permissions.

### See also:

[nvmlRestrictedAPI\\_t](#)

## `nvmlReturn_t nvmlDeviceSetFanSpeed_v2` (`nvmlDevice_t` device, unsigned int fan, unsigned int speed)

### Description

Sets the speed of a specified fan.

**WARNING:** This function changes the fan control policy to manual. It means that YOU have to monitor the temperature and adjust the fan speed accordingly. If you set the fan speed too low you can burn your GPU! Use `nvmlDeviceSetDefaultFanSpeed_v2` to restore default control policy.

For all cuda-capable discrete products with fans that are Maxwell or Newer.

device The identifier of the target device fan The index of the fan, starting at zero speed  
The target speed of the fan [0-100] in % of max speed

return `NVML_SUCCESS` if the fan speed has been set

`NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized

`NVML_ERROR_INVALID_ARGUMENT` if the device is not valid, or the speed

is outside acceptable ranges, or if the fan index doesn't reference an actual fan. `NVML_ERROR_NOT_SUPPORTED` if the device is older than Maxwell. `NVML_ERROR_UNKNOWN` if there was an unexpected error.

## `nvmlReturn_t nvmlDeviceSetGpcClkVfOffset` (`nvmlDevice_t device`, `int offset`)

### Parameters

#### `device`

The identifier of the target device

#### `offset`

The GPCCLK VF offset value to set

### Returns

- ▶ `NVML_SUCCESS` if offset has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or offset is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Set the GPCCLK VF offset value

## `nvmlReturn_t nvmlDeviceSetMemClkVfOffset` (`nvmlDevice_t device`, `int offset`)

### Parameters

#### `device`

The identifier of the target device

#### `offset`

The MemClk VF offset value to set

### Returns

- ▶ `NVML_SUCCESS` if offset has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or offset is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device does not support this feature

- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Set the MemClk (Memory Clock) VF offset value. It requires elevated privileges.

## `nvmlReturn_t` `nvmlDeviceSetConfComputeUnprotectedMemSize` (`nvmlDevice_t` device, unsigned long long sizeKiB)

### Parameters

#### `device`

Device Handle

#### `sizeKiB`

Unprotected Memory size to be set in KiB

### Returns

- ▶ `NVML_SUCCESS` if `sizeKiB` successfully set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device

### Description

Set Conf Computing Unprotected Memory Size.

For Ampere or newer fully supported devices. Supported on Linux, Windows TCC.

## `nvmlReturn_t` `nvmlSystemSetConfComputeGpusReadyState` (unsigned int isAcceptingWork)

### Parameters

#### `isAcceptingWork`

GPU accepting new work, `NVML_CC_ACCEPTING_CLIENT_REQUESTS_TRUE` or `NVML_CC_ACCEPTING_CLIENT_REQUESTS_FALSE`

**Description**

Set Conf Computing GPUs ready state.

For Ampere or newer fully supported devices. Supported on Linux, Windows TCC.

return

- ▶ `NVML_SUCCESS` if current GPUs ready state is successfully set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `isAcceptingWork` is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device

**nvmlReturn\_t**
**nvmlSystemSetConfComputeKeyRotationThresholdInfo  
(nvmlConfComputeSetKeyRotationThresholdInfo\_t  
\*pKeyRotationThrInfo)**
**Parameters****pKeyRotationThrInfo**

Reference to the key rotation threshold data

**Returns**

- ▶ `NVML_SUCCESS` if key rotation threshold max attacker advantage has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or memory is NULL
- ▶ `NVML_ERROR_INVALID_STATE` if confidential compute GPU ready state is enabled
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

**Description**

Set Conf Computing key rotation threshold.

HOPPER\_OR\_NEWER% Supported on Linux, Windows TCC.

This function is to set the confidential compute key rotation threshold parameters. `pKeyRotationThrInfo->maxAttackerAdvantage` should be in the range from `NVML_CC_KEY_ROTATION_THRESHOLD_ATTACKER_ADVANTAGE_MIN` to `NVML_CC_KEY_ROTATION_THRESHOLD_ATTACKER_ADVANTAGE_MAX`. Default value is 60.

## 4.18. NvLink Methods

This chapter describes methods that NVML can perform on NVLINK enabled devices.

**`nvmlReturn_t nvmlDeviceGetNvLinkState (nvmlDevice_t device, unsigned int link, nvmlEnableState_t *isActive)`**

### Parameters

#### **device**

The identifier of the target device

#### **link**

Specifies the NvLink link to be queried

#### **isActive**

`nvmlEnableState_t` where `NVML_FEATURE_ENABLED` indicates that the link is active and `NVML_FEATURE_DISABLED` indicates it is inactive

### Returns

- ▶ `NVML_SUCCESS` if `isActive` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device or link is invalid or `isActive` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device doesn't support this feature
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the state of the device's NvLink for the link specified

For Pascal or newer fully supported devices.

**`nvmlReturn_t nvmlDeviceGetNvLinkVersion (nvmlDevice_t device, unsigned int link, unsigned int *version)`**

### Parameters

#### **device**

The identifier of the target device

**link**

Specifies the NvLink link to be queried

**version**

Requested NvLink version

**Returns**

- ▶ NVML\_SUCCESS if version has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device or link is invalid or version is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieves the version of the device's NvLink for the link specified

For Pascal or newer fully supported devices.

**nvmlReturn\_t nvmlDeviceGetNvLinkCapability**  
**(nvmlDevice\_t device, unsigned int link,**  
**nvmlNvLinkCapability\_t capability, unsigned int**  
**\*capResult)**

**Parameters****device**

The identifier of the target device

**link**

Specifies the NvLink link to be queried

**capability**

Specifies the nvmlNvLinkCapability\_t to be queried

**capResult**

A boolean for the queried capability indicating that feature is available

**Returns**

- ▶ NVML\_SUCCESS if capResult has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, link, or capability is invalid or capResult is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature

- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the requested capability from the device's NvLink for the link specified. Please refer to the `nvmlNvLinkCapability_t` structure for the specific caps that can be queried. The return value should be treated as a boolean.

For Pascal or newer fully supported devices.

```
nvmlReturn_t nvmlDeviceGetNvLinkRemotePciInfo_v2
(nvmlDevice_t device, unsigned int link, nvmlPciInfo_t
*pci)
```

### Parameters

#### device

The identifier of the target device

#### link

Specifies the NvLink link to be queried

#### pci

`nvmlPciInfo_t` of the remote node for the specified link

### Returns

- ▶ NVML\_SUCCESS if pci has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device or link is invalid or pci is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves the PCI information for the remote node on a NvLink link. Note: `pciSubSystemId` is not filled in this function and is indeterminate.

For Pascal or newer fully supported devices.

```
nvmlReturn_t nvmlDeviceGetNvLinkErrorCounter
(nvmlDevice_t device, unsigned int link,
```



## `nvmlNvLinkErrorCounter_t` counter, unsigned long long \*counterValue)

### Parameters

#### **device**

The identifier of the target device

#### **link**

Specifies the NvLink link to be queried

#### **counter**

Specifies the NvLink counter to be queried

#### **counterValue**

Returned counter value

### Returns

- ▶ `NVML_SUCCESS` if counter has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device, link, or counter is invalid or counterValue is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the device doesn't support this feature
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the specified error counter value Please refer to `nvmlNvLinkErrorCounter_t` for error counters that are available

For Pascal or newer fully supported devices.

## `nvmlReturn_t` `nvmlDeviceResetNvLinkErrorCounters` (`nvmlDevice_t` device, unsigned int link)

### Parameters

#### **device**

The identifier of the target device

#### **link**

Specifies the NvLink link to be queried

### Returns

- ▶ `NVML_SUCCESS` if the reset is successful
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized

- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device or link is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Resets all error counters to zero Please refer to `nvmlNvLinkErrorCounter_t` for the list of error counters that are reset

For Pascal or newer fully supported devices.

`nvmlReturn_t nvmlDeviceSetNvLinkUtilizationControl`  
 (`nvmlDevice_t device`, `unsigned int link`, `unsigned int counter`, `nvmlNvLinkUtilizationControl_t *control`, `unsigned int reset`)

### Parameters

#### device

The identifier of the target device

#### link

Specifies the NvLink link to be queried

#### counter

Specifies the counter that should be set (0 or 1).

#### control

A reference to the `nvmlNvLinkUtilizationControl_t` to set

#### reset

Resets the counters on set if non-zero

### Returns

- ▶ NVML\_SUCCESS if the control has been set successfully
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, counter, link, or control is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Deprecated: Setting utilization counter control is no longer supported.

Set the NVLINK utilization counter control information for the specified counter, 0 or 1. Please refer to [nvmlNvLinkUtilizationControl\\_t](#) for the structure definition. Performs a reset of the counters if the reset parameter is non-zero.

For Pascal or newer fully supported devices.

**nvmlReturn\_t nvmlDeviceGetNvLinkUtilizationControl**  
(**nvmlDevice\_t** device, **unsigned int** link, **unsigned int** counter, **nvmlNvLinkUtilizationControl\_t \*control**)

### Parameters

#### device

The identifier of the target device

#### link

Specifies the NvLink link to be queried

#### counter

Specifies the counter that should be set (0 or 1).

#### control

A reference to the [nvmlNvLinkUtilizationControl\\_t](#) to place information

### Returns

- ▶ NVML\_SUCCESS if the control has been set successfully
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, counter, link, or control is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Deprecated: Getting utilization counter control is no longer supported.

Get the NVLINK utilization counter control information for the specified counter, 0 or 1. Please refer to [nvmlNvLinkUtilizationControl\\_t](#) for the structure definition

For Pascal or newer fully supported devices.

**nvmlReturn\_t nvmlDeviceGetNvLinkUtilizationCounter**  
(**nvmlDevice\_t** device, **unsigned int** link, **unsigned int**

counter, unsigned long long \*rxcounter, unsigned long long \*txcounter)

### Parameters

#### device

The identifier of the target device

#### link

Specifies the NvLink link to be queried

#### counter

Specifies the counter that should be read (0 or 1).

#### rxcounter

Receive counter return value

#### txcounter

Transmit counter return value

### Returns

- ▶ NVML\_SUCCESS if rxcounter and txcounter have been successfully set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, counter, or link is invalid or rxcounter or txcounter are NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Deprecated: Use [nvmlDeviceGetFieldValues](#) with NVML\_FI\_DEV\_NVLINK\_THROUGHPUT\_\* as field values instead.

Retrieve the NVLINK utilization counter based on the current control for a specified counter. In general it is good practice to use [nvmlDeviceSetNvLinkUtilizationControl](#) before reading the utilization counters as they have no default state

For Pascal or newer fully supported devices.

## `nvmlReturn_t` `nvmlDeviceFreezeNvLinkUtilizationCounter`

## `(nvmlDevice_t device, unsigned int link, unsigned int counter, nvmlEnableState_t freeze)`

### Parameters

#### **device**

The identifier of the target device

#### **link**

Specifies the NvLink link to be queried

#### **counter**

Specifies the counter that should be frozen (0 or 1).

#### **freeze**

NVML\_FEATURE\_ENABLED = freeze the receive and transmit counters

NVML\_FEATURE\_DISABLED = unfreeze the receive and transmit counters

### Returns

- ▶ NVML\_SUCCESS if counters were successfully frozen or unfrozen
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, link, counter, or freeze is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Deprecated: Freezing NVLINK utilization counters is no longer supported.

Freeze the NVLINK utilization counters Both the receive and transmit counters are operated on by this function

For Pascal or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceResetNvLinkUtilizationCounter (nvmlDevice_t device, unsigned int link, unsigned int counter)`

### Parameters

#### **device**

The identifier of the target device

#### **link**

Specifies the NvLink link to be reset

#### **counter**

Specifies the counter that should be reset (0 or 1)

**Returns**

- ▶ NVML\_SUCCESS if counters were successfully reset
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, link, or counter is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Deprecated: Resetting NVLINK utilization counters is no longer supported.

Reset the NVLINK utilization counters Both the receive and transmit counters are operated on by this function

For Pascal or newer fully supported devices.

```
nvmlReturn_t nvmlDeviceGetNvLinkRemoteDeviceType  
(nvmlDevice_t device, unsigned int link,  
nvmlIntNvLinkDeviceType_t *pNvLinkDeviceType)
```

**Parameters****device**

The device handle of the target GPU

**link**

The NVLink link index on the target GPU

**pNvLinkDeviceType**

Pointer in which the output remote device type is returned

**Returns**

- ▶ NVML\_SUCCESS if pNvLinkDeviceType has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if NVLink is not supported
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device or link is invalid, or pNvLinkDeviceType is NULL
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Get the NVLink device type of the remote device connected over the given link.

## 4.19. Event Handling Methods

This chapter describes methods that NVML can perform against each device to register and wait for some event to occur.

`struct nvmlEventData_t`

### Event Types

`typedef struct nvmlEventSet_st *nvmlEventSet_t`

Handle to an event set

`nvmlReturn_t nvmlEventSetCreate (nvmlEventSet_t *set)`

#### Parameters

`set`

Reference in which to return the event handle

#### Returns

- ▶ `NVML_SUCCESS` if the event has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if set is NULL
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

#### Description

Create an empty set of events. Event set should be freed by `nvmlEventSetFree`

For Fermi or newer fully supported devices.

#### See also:

`nvmlEventSetFree`

## `nvmlReturn_t nvmlDeviceRegisterEvents (nvmlDevice_t device, unsigned long long eventTypes, nvmlEventSet_t set)`

### Parameters

#### `device`

The identifier of the target device

#### `eventTypes`

Bitmask of [Event Types](#) to record

#### `set`

Set to which add new event types

### Returns

- ▶ `NVML_SUCCESS` if the event has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `eventTypes` is invalid or `set` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the platform does not support this feature or some of requested event types
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Starts recording of events on a specified devices and add the events to specified [nvmlEventSet\\_t](#)

For Fermi or newer fully supported devices. Ecc events are available only on ECC enabled devices (see [nvmlDeviceGetTotalEccErrors](#)) Power capping events are available only on Power Management enabled devices (see [nvmlDeviceGetPowerManagementMode](#))

For Linux only.

**IMPORTANT:** Operations on `set` are not thread safe

This call starts recording of events on specific device. All events that occurred before this call are not recorded. Checking if some event occurred can be done with [nvmlEventSetWait\\_v2](#)

If function reports `NVML_ERROR_UNKNOWN`, event set is in undefined state and should be freed. If function reports `NVML_ERROR_NOT_SUPPORTED`, event set can still be used. None of the requested `eventTypes` are registered in that case.



**See also:**[Event Types](#)[nvmlDeviceGetSupportedEventTypes](#)[nvmlEventSetWait](#)[nvmlEventSetFree](#)

## **`nvmlReturn_t nvmlDeviceGetSupportedEventTypes` (`nvmlDevice_t device`, `unsigned long long *eventTypes`)**

**Parameters****device**

The identifier of the target device

**eventTypes**

Reference in which to return bitmask of supported events

**Returns**

- ▶ `NVML_SUCCESS` if the `eventTypes` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `eventTypes` is `NULL`
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

**Description**

Returns information about events supported on device

For Fermi or newer fully supported devices.

Events are not supported on Windows. So this function returns an empty mask in `eventTypes` on Windows.

**See also:**[Event Types](#)[nvmlDeviceRegisterEvents](#)

## `nvmlReturn_t nvmlEventSetWait_v2 (nvmlEventSet_t set, nvmlEventData_t *data, unsigned int timeoutms)`

### Parameters

#### **set**

Reference to set of events to wait on

#### **data**

Reference in which to return event data

#### **timeoutms**

Maximum amount of wait time in milliseconds for registered event

### Returns

- ▶ `NVML_SUCCESS` if the data has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if data is `NULL`
- ▶ `NVML_ERROR_TIMEOUT` if no event arrived in specified timeout or interrupt arrived
- ▶ `NVML_ERROR_GPU_IS_LOST` if a GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Waits on events and delivers events

For Fermi or newer fully supported devices.

If some events are ready to be delivered at the time of the call, function returns immediately. If there are no events ready to be delivered, function sleeps till event arrives but not longer than specified timeout. This function in certain conditions can return before specified timeout passes (e.g. when interrupt arrives)

On Windows, in case of xid error, the function returns the most recent xid error type seen by the system. If there are multiple xid errors generated before `nvmlEventSetWait` is invoked then the last seen xid error type is returned for all xid error events.

On Linux, every xid error event would return the associated event data and other information if applicable.

In MIG mode, if device handle is provided, the API reports all the events for the available instances, only if the caller has appropriate privileges. In absence of required privileges, only the events which affect all the instances (i.e. whole device) are reported.

This API does not currently support per-instance event reporting using MIG device handles.

**See also:**

[Event Types](#)

[nvmlDeviceRegisterEvents](#)

## **nvmlReturn\_t nvmlEventSetFree (nvmlEventSet\_t set)**

### **Parameters**

**set**

Reference to events to be released

### **Returns**

- ▶ NVML\_SUCCESS if the event has been successfully released
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### **Description**

Releases events in the set

For Fermi or newer fully supported devices.

**See also:**

[nvmlDeviceRegisterEvents](#)

## **4.19.1. Event Types**

### Event Handling Methods

Event Types which user can be notified about. See description of particular functions for details.

See [nvmlDeviceRegisterEvents](#) and [nvmlDeviceGetSupportedEventTypes](#) to check which devices support each event.

Types can be combined with bitwise or operator '|' when passed to [nvmlDeviceRegisterEvents](#)

```
#define nvmlEventTypeSingleBitEccError 0x0000000000000001LL
```

Event about single bit ECC errors.



A corrected texture memory error is not an ECC error, so it does not generate a single bit event

```
#define nvmlEventTypeDoubleBitEccError 0x0000000000000002LL
```

Event about double bit ECC errors.



An uncorrected texture memory error is not an ECC error, so it does not generate a double bit event

```
#define nvmlEventTypePState 0x0000000000000004LL
```

Event about PState changes.



On Fermi architecture PState changes are also an indicator that GPU is throttling down due to no work being executed on the GPU, power capping or thermal capping. In a typical situation, Fermi-based GPU should stay in P0 for the duration of the execution of the compute process.

```
#define nvmlEventTypeXidCriticalError 0x0000000000000008LL
```

Event that Xid critical error occurred.

```
#define nvmlEventTypeClock 0x0000000000000010LL
```

Event about clock changes.

Kepler only

```
#define nvmlEventTypePowerSourceChange 0x0000000000000080LL
```

Event about AC/Battery power source changes.

```
#define nvmlEventMigConfigChange 0x0000000000000100LL
```

Event about MIG configuration changes.

```
#define nvmlEventTypeNone 0x0000000000000000LL
```

Mask with no events.

```
#define nvmlEventTypeAll (nvmlEventTypeNone
\ | nvmlEventTypeSingleBitEccError \ |
nvmlEventTypeDoubleBitEccError \ | nvmlEventTypePState \
| nvmlEventTypeClock \ | nvmlEventTypeXidCriticalError \ |
nvmlEventTypePowerSourceChange \ | nvmlEventMigConfigChange
\ )
```

Mask of all events.

## 4.20. Drain states

This chapter describes methods that NVML can perform against each device to control their drain state and recognition by NVML and NVIDIA kernel driver. These methods can be used with out-of-band tools to power on/off GPUs, enable robust reset scenarios, etc.

```
nvmlReturn_t nvmlDeviceModifyDrainState
(nvmlPciInfo_t *pciInfo, nvmlEnableState_t newState)
```

### Parameters

#### pciInfo

The PCI address of the GPU drain state to be modified

#### newState

The drain state that should be entered, see [nvmlEnableState\\_t](#)

### Returns

- ▶ NVML\_SUCCESS if counters were successfully reset
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if nvmlIndex or newState is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_NO\_PERMISSION if the calling process has insufficient permissions to perform operation
- ▶ NVML\_ERROR\_IN\_USE if the device has persistence mode turned on
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Modify the drain state of a GPU. This method forces a GPU to no longer accept new incoming requests. Any new NVML process will no longer see this GPU. Persistence

mode for this GPU must be turned off before this call is made. Must be called as administrator. For Linux only.

For Pascal or newer fully supported devices. Some Kepler devices supported.

```
nvmlReturn_t nvmlDeviceQueryDrainState  
(nvmlPciInfo_t *pciInfo, nvmlEnableState_t  
*currentState)
```

### Parameters

#### **pciInfo**

The PCI address of the GPU drain state to be queried

#### **currentState**

The current drain state for this GPU, see [nvmlEnableState\\_t](#)

### Returns

- ▶ NVML\_SUCCESS if counters were successfully reset
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if nvmlIndex or currentState is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Query the drain state of a GPU. This method is used to check if a GPU is in a currently draining state. For Linux only.

For Pascal or newer fully supported devices. Some Kepler devices supported.

```
nvmlReturn_t nvmlDeviceRemoveGpu_v2 (nvmlPciInfo_t  
*pciInfo, nvmlDetachGpuState_t gpuState,  
nvmlPcieLinkState_t linkState)
```

### Parameters

#### **pciInfo**

The PCI address of the GPU to be removed

#### **gpuState**

Whether the GPU is to be removed, from the OS see [nvmlDetachGpuState\\_t](#)

#### **linkState**

Requested upstream PCIe link state, see [nvmlPcieLinkState\\_t](#)

## Returns

- ▶ NVML\_SUCCESS if counters were successfully reset
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if nvmlIndex is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the device doesn't support this feature
- ▶ NVML\_ERROR\_IN\_USE if the device is still in use and cannot be removed

## Description

This method will remove the specified GPU from the view of both NVML and the NVIDIA kernel driver as long as no other processes are attached. If other processes are attached, this call will return NVML\_ERROR\_IN\_USE and the GPU will be returned to its original "draining" state. Note: the only situation where a process can still be attached after `nvmlDeviceModifyDrainState()` is called to initiate the draining state is if that process was using, and is still using, a GPU before the call was made. Also note, persistence mode counts as an attachment to the GPU thus it must be disabled prior to this call.

For long-running NVML processes please note that this will change the enumeration of current GPUs. For example, if there are four GPUs present and GPU1 is removed, the new enumeration will be 0-2. Also, device handles after the removed GPU will not be valid and must be re-established. Must be run as administrator. For Linux only.

For Pascal or newer fully supported devices. Some Kepler devices supported.

## `nvmlReturn_t nvmlDeviceDiscoverGpus (nvmlPciInfo_t *pciInfo)`

### Parameters

#### `pciInfo`

The PCI tree to be searched. Only the domain, bus, and device fields are used in this call.

### Returns

- ▶ NVML\_SUCCESS if counters were successfully reset
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if pciInfo is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if the operating system does not support this feature
- ▶ NVML\_ERROR\_OPERATING\_SYSTEM if the operating system is denying this feature

- ▶ `NVML_ERROR_NO_PERMISSION` if the calling process has insufficient permissions to perform operation
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Request the OS and the NVIDIA kernel driver to rediscover a portion of the PCI subsystem looking for GPUs that were previously removed. The portion of the PCI tree can be narrowed by specifying a domain, bus, and device. If all are zeroes then the entire PCI tree will be searched. Please note that for long-running NVML processes the enumeration will change based on how many GPUs are discovered and where they are inserted in bus order.

In addition, all newly discovered GPUs will be initialized and their ECC scrubbed which may take several seconds per GPU. Also, all device handles are no longer guaranteed to be valid post discovery.

Must be run as administrator. For Linux only.

For Pascal or newer fully supported devices. Some Kepler devices supported.

## 4.21. Field Value Queries

This chapter describes NVML operations that are associated with retrieving Field Values from NVML

**`nvmlReturn_t nvmlDeviceGetFieldValues (nvmlDevice_t device, int valuesCount, nvmlFieldValue_t *values)`**

### Parameters

#### **device**

The device handle of the GPU to request field values for

#### **valuesCount**

Number of entries in values that should be retrieved

#### **values**

Array of valuesCount structures to hold field values. Each value's fieldId must be populated prior to this call

### Returns

- ▶ `NVML_SUCCESS` if any values in values were populated. Note that you must check the `nvmlReturn` field of each value for each individual status
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or values is NULL



**Description**

Request values for a list of fields for a device. This API allows multiple fields to be queried at once. If any of the underlying fieldIds are populated by the same driver call, the results for those field IDs will be populated from a single call rather than making a driver call for each fieldId.

```
nvmlReturn_t nvmlDeviceClearFieldValues
(nvmlDevice_t device, int valuesCount,
nvmlFieldValue_t *values)
```

**Parameters****device**

The device handle of the GPU to request field values for

**valuesCount**

Number of entries in values that should be cleared

**values**

Array of valuesCount structures to hold field values. Each value's fieldId must be populated prior to this call

**Returns**

- ▶ NVML\_SUCCESS if any values in values were cleared. Note that you must check the nvmlReturn field of each value for each individual status
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or values is NULL

**Description**

Clear values for a list of fields for a device. This API allows multiple fields to be cleared at once.

## 4.22. Enums, Constants and Structs

## 4.23. vGPU APIs

This chapter describes operations that are associated with NVIDIA vGPU Software products.

## `nvmlReturn_t nvmlDeviceGetVirtualizationMode` (`nvmlDevice_t device`, `nvmlGpuVirtualizationMode_t *pVirtualMode`)

### Parameters

#### `device`

Identifier of the target device

#### `pVirtualMode`

Reference to virtualization mode. One of `NVML_GPU_VIRTUALIZATION_?`

### Returns

- ▶ `NVML_SUCCESS` if `pVirtualMode` is fetched
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `device` is invalid or `pVirtualMode` is `NULL`
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

This method is used to get the virtualization mode corresponding to the GPU.

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetHostVgpuMode` (`nvmlDevice_t device`, `nvmlHostVgpuMode_t *pHostVgpuMode`)

### Parameters

#### `device`

The identifier of the target device

#### `pHostVgpuMode`

Reference in which to return the current vGPU mode

### Returns

- ▶ `NVML_SUCCESS` if device's vGPU mode has been successfully retrieved
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device handle is 0 or `pVgpuMode` is `NULL`

- ▶ `NVML_ERROR_NOT_SUPPORTED` if device doesn't support this feature.
- ▶ `NVML_ERROR_UNKNOWN` if any unexpected error occurred

### Description

Queries if SR-IOV host operation is supported on a vGPU supported device.

Checks whether SR-IOV host capability is supported by the device and the driver, and indicates device is in SR-IOV mode if both of these conditions are true.

## `nvmlReturn_t nvmlDeviceSetVirtualizationMode` (`nvmlDevice_t device`, `nvmlGpuVirtualizationMode_t virtualMode`)

### Parameters

#### `device`

Identifier of the target device

#### `virtualMode`

virtualization mode. One of `NVML_GPU_VIRTUALIZATION_?`

### Returns

- ▶ `NVML_SUCCESS` if `virtualMode` is set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid or `virtualMode` is `NULL`
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_NOT_SUPPORTED` if setting of virtualization mode is not supported.
- ▶ `NVML_ERROR_NO_PERMISSION` if setting of virtualization mode is not allowed for this client.

### Description

This method is used to set the virtualization mode corresponding to the GPU.

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetVgpuHeterogeneousMode(nvmlDevice_t device, nvmlVgpuHeterogeneousMode_t *pHeterogeneousMode)`

### Parameters

#### **device**

The identifier of the target device

#### **pHeterogeneousMode**

Pointer to the caller-provided structure of `nvmlVgpuHeterogeneousMode_t`

### Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device is invalid or `pHeterogeneousMode` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't support this feature
- ▶ `NVML_ERROR_VERSION_MISMATCH` If the version of `pHeterogeneousMode` is invalid
- ▶ `NVML_ERROR_UNKNOWN` On any unexpected error

### Description

Get the vGPU heterogeneous mode for the device.

When in heterogeneous mode, a vGPU can concurrently host timesliced vGPUs with differing framebuffer sizes.

On successful return, the function returns `pHeterogeneousMode->mode` with the current vGPU heterogeneous mode. `pHeterogeneousMode->version` is the version number of the structure `nvmlVgpuHeterogeneousMode_t`, the caller should set the correct version number to retrieve the vGPU heterogeneous mode. `pHeterogeneousMode->mode` can either be `NVML_FEATURE_ENABLED` or `NVML_FEATURE_DISABLED`.

## `nvmlReturn_t nvmlDeviceSetVgpuHeterogeneousMode (nvmlDevice_t device, const nvmlVgpuHeterogeneousMode_t *pHeterogeneousMode)`

### Parameters

#### **device**

Identifier of the target device

#### **pHeterogeneousMode**

Pointer to the caller-provided structure of `nvmlVgpuHeterogeneousMode_t`

### Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device or `pHeterogeneousMode` is `NULL` or `pHeterogeneousMode->mode` is invalid
- ▶ `NVML_ERROR_IN_USE` If the device is in use
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation
- ▶ `NVML_ERROR_NOT_SUPPORTED` If MIG is enabled or device doesn't support this feature
- ▶ `NVML_ERROR_VERSION_MISMATCH` If the version of `pHeterogeneousMode` is invalid
- ▶ `NVML_ERROR_UNKNOWN` On any unexpected error

### Description

Enable or disable vGPU heterogeneous mode for the device.

When in heterogeneous mode, a vGPU can concurrently host timesliced vGPUs with differing framebuffer sizes.

API would return an appropriate error code upon unsuccessful activation. For example, the heterogeneous mode set will fail with error `NVML_ERROR_IN_USE` if any vGPU instance is active on the device. The caller of this API is expected to shutdown the vGPU VMs and retry setting the mode. On successful return, the function updates the vGPU heterogeneous mode with the user provided `pHeterogeneousMode->mode`. `pHeterogeneousMode->version` is the version number of the structure `nvmlVgpuHeterogeneousMode_t`, the caller should set the correct version number to set the vGPU heterogeneous mode.

## nvmlReturn\_t nvmlVgpuInstanceGetPlacementId (nvmlVgpuInstance\_t vgpuInstance, nvmlVgpuPlacementId\_t \*pPlacement)

### Parameters

#### vgpuInstance

Identifier of the target vGPU instance

#### pPlacement

Pointer to vGPU placement ID structure nvmlVgpuPlacementId\_t

### Returns

- ▶ NVML\_SUCCESS If information is successfully retrieved
- ▶ NVML\_ERROR\_NOT\_FOUND If vgpuInstance does not match a valid active vGPU instance
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT If vgpuInstance is invalid or pPlacement is NULL
- ▶ NVML\_ERROR\_VERSION\_MISMATCH If the version of pPlacement is invalid
- ▶ NVML\_ERROR\_UNKNOWN On any unexpected error

### Description

Query the placement ID of active vGPU instance.

When in vGPU heterogeneous mode, this function returns a valid placement ID as pPlacement->placementId else NVML\_INVALID\_VGPU\_PLACEMENT\_ID is returned. pPlacement->version is the version number of the structure nvmlVgpuPlacementId\_t, the caller should set the correct version number to get placement id of the vGPU instance vgpuInstance.

## nvmlReturn\_t nvmlDeviceGetVgpuTypeSupportedPlacements (nvmlDevice\_t device, nvmlVgpuTypeId\_t vgpuTypeId, nvmlVgpuPlacementList\_t \*pPlacementList)

### Parameters

#### device

Identifier of the target device

#### vgpuTypeId

Handle to vGPU type. The vGPU type ID

**pPlacementList**

Pointer to the vGPU placement structure `nvmlVgpuPlacementList_t`

**Returns**

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device or `vgpuTypeId` is invalid or `pPlacementList` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device or `vgpuTypeId` isn't supported
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation
- ▶ `NVML_ERROR_VERSION_MISMATCH` If the version of `pPlacementList` is invalid
- ▶ `NVML_ERROR_UNKNOWN` On any unexpected error

**Description**

Query the supported vGPU placement ID of the vGPU type.

An array of supported vGPU placement IDs for the vGPU type ID indicated by `vgpuTypeId` is returned in the caller-supplied buffer of `pPlacementList->placementIds`. Memory needed for the `placementIds` array should be allocated based on maximum instances of a vGPU type which can be queried via `nvmlVgpuTypeGetMaxInstances()`.

This function will return supported placement IDs even if GPU is not in vGPU heterogeneous mode.

**nvmlReturn\_t****nvmlDeviceGetVgpuTypeCreatablePlacements**

(`nvmlDevice_t` device, `nvmlVgpuTypeId_t` vgpuTypeId, `nvmlVgpuPlacementList_t` \*pPlacementList)

**Parameters****device**

The identifier of the target device

**vgpuTypeId**

Handle to vGPU type. The vGPU type ID

**pPlacementList**

Pointer to the list of vGPU placement structure `nvmlVgpuPlacementList_t`

**Returns**

- ▶ `NVML_SUCCESS` Upon success

- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device or `vgpuTypeId` is invalid or `pPlacementList` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device or `vgpuTypeId` isn't supported
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation
- ▶ `NVML_ERROR_VERSION_MISMATCH` If the version of `pPlacementList` is invalid
- ▶ `NVML_ERROR_UNKNOWN` On any unexpected error

### Description

Query the creatable vGPU placement ID of the vGPU type.

An array of creatable vGPU placement IDs for the vGPU type ID indicated by `vgpuTypeId` is returned in the caller-supplied buffer of `pPlacementList->placementIds`. Memory needed for the `placementIds` array should be allocated based on maximum instances of a vGPU type which can be queried via `nvmlVgpuTypeGetMaxInstances()`. The creatable vGPU placement IDs may differ over time, as there may be restrictions on what type of vGPU the vGPU instance is running.

The function will return `NVML_ERROR_NOT_SUPPORTED` if the device is not in vGPU heterogeneous mode.

## `nvmlReturn_t nvmlVgpuTypeGetGspHeapSize` (`nvmlVgpuTypeId_t vgpuTypeId`, unsigned long long `*gspHeapSize`)

### Parameters

#### `vgpuTypeId`

Handle to vGPU type

#### `gspHeapSize`

Reference to return the GSP heap size value

### Returns

- ▶ `NVML_SUCCESS` Successful completion
- ▶ `NVML_ERROR_UNINITIALIZED` If the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If `vgpuTypeId` is invalid, or `gspHeapSize` is `NULL`
- ▶ `NVML_ERROR_UNKNOWN` On any unexpected error



**Description**

Retrieve the static GSP heap size of the vGPU type in bytes

**nvmlReturn\_t nvmlVgpuTypeGetFbReservation**  
(**nvmlVgpuTypeId\_t vgpuTypeId**, unsigned long long **\*fbReservation**)

**Parameters****vgpuTypeId**

Handle to vGPU type

**fbReservation**

Reference to return the framebuffer reservation

**Returns**

- ▶ **NVML\_SUCCESS** Successful completion
- ▶ **NVML\_ERROR\_UNINITIALIZED** If the library has not been successfully initialized
- ▶ **NVML\_ERROR\_INVALID\_ARGUMENT** If **vgpuTypeId** is invalid, or **fbReservation** is NULL
- ▶ **NVML\_ERROR\_UNKNOWN** On any unexpected error

**Description**

Retrieve the static framebuffer reservation of the vGPU type in bytes

**nvmlReturn\_t nvmlDeviceSetVgpuCapabilities**  
(**nvmlDevice\_t device**, **nvmlDeviceVgpuCapability\_t capability**, **nvmlEnableState\_t state**)

**Parameters****device**

The identifier of the target device

**capability**

Specifies the **nvmlDeviceVgpuCapability\_t** to be set

**state**

The target capability mode

**Returns**

- ▶ **NVML\_SUCCESS** Successful completion

- ▶ `NVML_ERROR_UNINITIALIZED` If the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device is invalid, or capability is invalid, or state is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` The API is not supported in current state, or device not in vGPU mode
- ▶ `NVML_ERROR_UNKNOWN` On any unexpected error

### Description

Set the desirable vGPU capability of a device

Refer to the `nvmlDeviceVgpuCapability_t` structure for the specific capabilities that can be set. See `nvmlEnableState_t` for available state.

## `nvmlReturn_t nvmlDeviceGetGridLicensableFeatures_v4` (`nvmlDevice_t device, nvmlGridLicensableFeatures_t *pGridLicensableFeatures`)

### Parameters

#### `device`

Identifier of the target device

#### `pGridLicensableFeatures`

Pointer to structure in which vGPU software licensable features are returned

### Returns

- ▶ `NVML_SUCCESS` if licensable features are successfully retrieved
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `pGridLicensableFeatures` is `NULL`
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the vGPU Software licensable features.

Identifies whether the system supports vGPU Software Licensing. If it does, return the list of licensable feature(s) and their current license status.

## 4.24. vGPU Management

This chapter describes APIs supporting NVIDIA vGPU.

## `nvmlReturn_t nvmlGetVgpuDriverCapabilities` (`nvmlVgpuDriverCapability_t` `capability`, unsigned int `*capResult`)

### Parameters

#### `capability`

Specifies the `nvmlVgpuDriverCapability_t` to be queried

#### `capResult`

A boolean for the queried capability indicating that feature is supported

### Returns

- ▶ `NVML_SUCCESS` successful completion
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `capability` is invalid, or `capResult` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` the API is not supported in current state or devices not in vGPU mode
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the requested vGPU driver capability.

Refer to the `nvmlVgpuDriverCapability_t` structure for the specific capabilities that can be queried. The return value in `capResult` should be treated as a boolean, with a non-zero value indicating that the capability is supported.

For Maxwell or newer fully supported devices.

## `nvmlReturn_t nvmlDeviceGetVgpuCapabilities` (`nvmlDevice_t` `device`, `nvmlDeviceVgpuCapability_t` `capability`, unsigned int `*capResult`)

### Parameters

#### `device`

The identifier of the target device

#### `capability`

Specifies the `nvmlDeviceVgpuCapability_t` to be queried

#### `capResult`

Specifies that the queried capability is supported, and also returns capability's data

**Returns**

- ▶ NVML\_SUCCESS successful completion
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, or capability is invalid, or capResult is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED the API is not supported in current state or device not in vGPU mode
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve the requested vGPU capability for GPU.

Refer to the `nvmlDeviceVgpuCapability_t` structure for the specific capabilities that can be queried. The return value in `capResult` reports a non-zero value indicating that the capability is supported, and also reports the capability's data based on the queried capability.

For Maxwell or newer fully supported devices.

**`nvmlReturn_t nvmlDeviceGetSupportedVgpus`**  
**(`nvmlDevice_t device`, `unsigned int *vgpuCount`,**  
**`nvmlVgpuTypeId_t *vgpuTypeIds`)**

**Parameters****device**

The identifier of the target device

**vgpuCount**

Pointer to caller-supplied array size, and returns number of vGPU types

**vgpuTypeIds**

Pointer to caller-supplied array in which to return list of vGPU types

**Returns**

- ▶ NVML\_SUCCESS successful completion
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE `vgpuTypeIds` buffer is too small, array element count is returned in `vgpuCount`
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if `vgpuCount` is NULL or device is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if vGPU is not supported by the device
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

## Description

Retrieve the supported vGPU types on a physical GPU (device).

An array of supported vGPU types for the physical GPU indicated by device is returned in the caller-supplied buffer pointed at by `vgpuTypeIds`. The element count of `nvmlVgpuTypeId_t` array is passed in `vgpuCount`, and `vgpuCount` is used to return the number of vGPU types written to the buffer.

If the supplied buffer is not large enough to accommodate the vGPU type array, the function returns `NVML_ERROR_INSUFFICIENT_SIZE`, with the element count of `nvmlVgpuTypeId_t` array required in `vgpuCount`. To query the number of vGPU types supported for the GPU, call this function with `*vgpuCount = 0`. The code will return `NVML_ERROR_INSUFFICIENT_SIZE`, or `NVML_SUCCESS` if no vGPU types are supported.

```
nvmlReturn_t nvmlDeviceGetCreatableVgpus
(nvmlDevice_t device, unsigned int *vgpuCount,
nvmlVgpuTypeId_t *vgpuTypeIds)
```

## Parameters

### **device**

The identifier of the target device

### **vgpuCount**

Pointer to caller-supplied array size, and returns number of vGPU types

### **vgpuTypeIds**

Pointer to caller-supplied array in which to return list of vGPU types

## Returns

- ▶ `NVML_SUCCESS` successful completion
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` `vgpuTypeIds` buffer is too small, array element count is returned in `vgpuCount`
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuCount` is `NULL`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if vGPU is not supported by the device
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

## Description

Retrieve the currently creatable vGPU types on a physical GPU (device).

An array of creatable vGPU types for the physical GPU indicated by device is returned in the caller-supplied buffer pointed at by `vgpuTypeIds`. The element count of

`nvmlVgpuTypeId_t` array is passed in `vgpuCount`, and `vgpuCount` is used to return the number of vGPU types written to the buffer.

The creatable vGPU types for a device may differ over time, as there may be restrictions on what type of vGPU types can concurrently run on a device. For example, if only one vGPU type is allowed at a time on a device, then the creatable list will be restricted to whatever vGPU type is already running on the device.

If the supplied buffer is not large enough to accommodate the vGPU type array, the function returns `NVML_ERROR_INSUFFICIENT_SIZE`, with the element count of `nvmlVgpuTypeId_t` array required in `vgpuCount`. To query the number of vGPU types that can be created for the GPU, call this function with `*vgpuCount = 0`. The code will return `NVML_ERROR_INSUFFICIENT_SIZE`, or `NVML_SUCCESS` if no vGPU types are creatable.

## `nvmlReturn_t nvmlVgpuTypeGetClass` (`nvmlVgpuTypeId_t vgpuTypeId`, `char *vgpuTypeClass`, `unsigned int *size`)

### Parameters

#### `vgpuTypeId`

Handle to vGPU type

#### `vgpuTypeClass`

Pointer to string array to return class in

#### `size`

Size of string

### Returns

- ▶ `NVML_SUCCESS` successful completion
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuTypeId` is invalid, or `vgpuTypeClass` is `NULL`
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `size` is too small
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the class of a vGPU type. It will not exceed 64 characters in length (including the NUL terminator). See [nvmlConstants::NVML\\_DEVICE\\_NAME\\_BUFFER\\_SIZE](#).

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuTypeGetName` (`nvmlVgpuTypeId_t vgpuTypeId`, `char *vgpuTypeName`, `unsigned int *size`)

### Parameters

#### `vgpuTypeId`

Handle to vGPU type

#### `vgpuTypeName`

Pointer to buffer to return name

#### `size`

Size of buffer

### Returns

- ▶ `NVML_SUCCESS` successful completion
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuTypeId` is invalid, or name is `NULL`
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `size` is too small
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the vGPU type name.

The name is an alphanumeric string that denotes a particular vGPU, e.g. GRID M60-2Q. It will not exceed 64 characters in length (including the NUL terminator). See [nvmlConstants::NVML\\_DEVICE\\_NAME\\_BUFFER\\_SIZE](#).

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuTypeGetGpuInstanceId` (`nvmlVgpuTypeId_t vgpuTypeId`, `unsigned int` `*gpuInstanceId`)

### Parameters

#### `vgpuTypeId`

Handle to vGPU type

#### `gpuInstanceId`

GPU Instance Profile ID

**Returns**

- ▶ NVML\_SUCCESS successful completion
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if device is not in vGPU Host virtualization mode
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if vgpuTypeId is invalid, or gpuInstanceId is NULL
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve the GPU Instance Profile ID for the given vGPU type ID. The API will return a valid GPU Instance Profile ID for the MIG capable vGPU types, else INVALID\_GPU\_INSTANCE\_PROFILE\_ID is returned.

For Kepler or newer fully supported devices.

**nvmlReturn\_t nvmlVgpuTypeGetDeviceID**  
**(nvmlVgpuTypeId\_t vgpuTypeId, unsigned long long**  
**\*deviceID, unsigned long long \*subsystemID)**

**Parameters****vgpuTypeId**

Handle to vGPU type

**deviceID**

Device ID and vendor ID of the device contained in single 32 bit value

**subsystemID**

Subsystem ID and subsystem vendor ID of the device contained in single 32 bit value

**Returns**

- ▶ NVML\_SUCCESS successful completion
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if vgpuTypeId is invalid, or deviceID or subsystemID are NULL
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve the device ID of a vGPU type.

For Kepler or newer fully supported devices.



## `nvmlReturn_t nvmlVgpuTypeGetFramebufferSize` (`nvmlVgpuTypeId_t vgpuTypeId`, unsigned long long `*fbSize`)

### Parameters

#### `vgpuTypeId`

Handle to vGPU type

#### `fbSize`

Pointer to framebuffer size in bytes

### Returns

- ▶ `NVML_SUCCESS` successful completion
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuTypeId` is invalid, or `fbSize` is `NULL`
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the vGPU framebuffer size in bytes.

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuTypeGetNumDisplayHeads` (`nvmlVgpuTypeId_t vgpuTypeId`, unsigned int `*numDisplayHeads`)

### Parameters

#### `vgpuTypeId`

Handle to vGPU type

#### `numDisplayHeads`

Pointer to number of display heads

### Returns

- ▶ `NVML_SUCCESS` successful completion
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuTypeId` is invalid, or `numDisplayHeads` is `NULL`

- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieve count of vGPU's supported display heads.

For Kepler or newer fully supported devices.

**nvmlReturn\_t nvmlVgpuTypeGetResolution**  
**(nvmlVgpuTypeId\_t vgpuTypeId, unsigned int**  
**displayIndex, unsigned int \*xdim, unsigned int \*ydim)**

### Parameters

#### **vgpuTypeId**

Handle to vGPU type

#### **displayIndex**

Zero-based index of display head

#### **xdim**

Pointer to maximum number of pixels in X dimension

#### **ydim**

Pointer to maximum number of pixels in Y dimension

### Returns

- ▶ NVML\_SUCCESS successful completion
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if vgpuTypeId is invalid, or xdim or ydim are NULL, or displayIndex is out of range.
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieve vGPU display head's maximum supported resolution.

For Kepler or newer fully supported devices.

**nvmlReturn\_t nvmlVgpuTypeGetLicense**  
 (nvmlVgpuTypeId\_t vgpuTypeId, char  
 \*vgpuTypeLicenseString, unsigned int size)

### Parameters

#### **vgpuTypeId**

Handle to vGPU type

#### **vgpuTypeLicenseString**

Pointer to buffer to return license info

#### **size**

Size of vgpuTypeLicenseString buffer

### Returns

- ▶ NVML\_SUCCESS successful completion
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if vgpuTypeId is invalid, or vgpuTypeLicenseString is NULL
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if size is too small
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieve license requirements for a vGPU type

The license type and version required to run the specified vGPU type is returned as an alphanumeric string, in the form "<license name>,<version>", for example "GRID-Virtual-PC,2.0". If a vGPU is runnable with\* more than one type of license, the licenses are delimited by a semicolon, for example "GRID-Virtual-PC,2.0;GRID-Virtual-WS,2.0;GRID-Virtual-WS-Ext,2.0".

The total length of the returned string will not exceed 128 characters, including the NUL terminator. See [nvmlVgpuConstants::NVML\\_GRID\\_LICENSE\\_BUFFER\\_SIZE](#).

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuTypeGetFrameRateLimit` (`nvmlVgpuTypeId_t vgpuTypeId`, `unsigned int *frameRateLimit`)

### Parameters

#### `vgpuTypeId`

Handle to vGPU type

#### `frameRateLimit`

Reference to return the frame rate limit value

### Returns

- ▶ `NVML_SUCCESS` successful completion
- ▶ `NVML_ERROR_NOT_SUPPORTED` if frame rate limiter is turned off for the vGPU type
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuTypeId` is invalid, or `frameRateLimit` is `NULL`
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the static frame rate limit value of the vGPU type

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuTypeGetMaxInstances` (`nvmlDevice_t device`, `nvmlVgpuTypeId_t vgpuTypeId`, `unsigned int *vgpuInstanceCount`)

### Parameters

#### `device`

The identifier of the target device

#### `vgpuTypeId`

Handle to vGPU type

#### `vgpuInstanceCount`

Pointer to get the max number of vGPU instances that can be created on a device for given `vgpuTypeId`

**Returns**

- ▶ NVML\_SUCCESS successful completion
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if vgpuTypeId is invalid or is not supported on target device, or vgpuInstanceCount is NULL
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve the maximum number of vGPU instances creatable on a device for given vGPU type

For Kepler or newer fully supported devices.

**nvmlReturn\_t nvmlVgpuTypeGetMaxInstancesPerVm  
(nvmlVgpuTypeId\_t vgpuTypeId, unsigned int  
\*vgpuInstanceCountPerVm)**

**Parameters****vgpuTypeId**

Handle to vGPU type

**vgpuInstanceCountPerVm**

Pointer to get the max number of vGPU instances supported per VM for given vgpuTypeId

**Returns**

- ▶ NVML\_SUCCESS successful completion
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if vgpuTypeId is invalid, or vgpuInstanceCountPerVm is NULL
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve the maximum number of vGPU instances supported per VM for given vGPU type

For Kepler or newer fully supported devices.

**nvmlReturn\_t nvmlDeviceGetActiveVgpus (nvmlDevice\_t device, unsigned int \*vgpuCount, nvmlVgpuInstance\_t \*vgpuInstances)**

### Parameters

#### device

The identifier of the target device

#### vgpuCount

Pointer which passes in the array size as well as get back the number of types

#### vgpuInstances

Pointer to array in which to return list of vGPU instances

### Returns

- ▶ NVML\_SUCCESS successful completion
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, or vgpuCount is NULL
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if size is too small
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if vGPU is not supported by the device
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieve the active vGPU instances on a device.

An array of active vGPU instances is returned in the caller-supplied buffer pointed at by vgpuInstances. The array element count is passed in vgpuCount, and vgpuCount is used to return the number of vGPU instances written to the buffer.

If the supplied buffer is not large enough to accommodate the vGPU instance array, the function returns NVML\_ERROR\_INSUFFICIENT\_SIZE, with the element count of nvmlVgpuInstance\_t array required in vgpuCount. To query the number of active vGPU instances, call this function with \*vgpuCount = 0. The code will return NVML\_ERROR\_INSUFFICIENT\_SIZE, or NVML\_SUCCESS if no vGPU Types are supported.

For Kepler or newer fully supported devices.

**nvmlReturn\_t nvmlVgpuInstanceGetVmID**  
(**nvmlVgpuInstance\_t vgpuInstance, char \*vmId,**  
**unsigned int size, nvmlVgpuVmIdType\_t \*vmIdType**)

### Parameters

#### **vgpuInstance**

Identifier of the target vGPU instance

#### **vmId**

Pointer to caller-supplied buffer to hold VM ID

#### **size**

Size of buffer in bytes

#### **vmIdType**

Pointer to hold VM ID type

### Returns

- ▶ **NVML\_SUCCESS** successful completion
- ▶ **NVML\_ERROR\_UNINITIALIZED** if the library has not been successfully initialized
- ▶ **NVML\_ERROR\_INVALID\_ARGUMENT** if **vmId** or **vmIdType** is **NULL**, or **vgpuInstance** is 0
- ▶ **NVML\_ERROR\_NOT\_FOUND** if **vgpuInstance** does not match a valid active vGPU instance on the system
- ▶ **NVML\_ERROR\_INSUFFICIENT\_SIZE** if **size** is too small
- ▶ **NVML\_ERROR\_UNKNOWN** on any unexpected error

### Description

Retrieve the VM ID associated with a vGPU instance.

The VM ID is returned as a string, not exceeding 80 characters in length (including the NUL terminator). See [nvmlConstants::NVML\\_DEVICE\\_UUID\\_BUFFER\\_SIZE](#).

The format of the VM ID varies by platform, and is indicated by the type identifier returned in **vmIdType**.

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuInstanceGetUUID` (`nvmlVgpuInstance_t vgpuInstance`, `char *uuid`, `unsigned int size`)

### Parameters

#### `vgpuInstance`

Identifier of the target vGPU instance

#### `uuid`

Pointer to caller-supplied buffer to hold vGPU UUID

#### `size`

Size of buffer in bytes

### Returns

- ▶ `NVML_SUCCESS` successful completion
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `uuid` is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `size` is too small
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the UUID of a vGPU instance.

The UUID is a globally unique identifier associated with the vGPU, and is returned as a 5-part hexadecimal string, not exceeding 80 characters in length (including the `NULL` terminator). See `nvmlConstants::NVML_DEVICE_UUID_BUFFER_SIZE`.

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuInstanceGetVmDriverVersion` (`nvmlVgpuInstance_t vgpuInstance`, `char *version`, `unsigned int length`)

### Parameters

#### `vgpuInstance`

Identifier of the target vGPU instance



**version**

Caller-supplied buffer to return driver version string

**length**

Size of version buffer

**Returns**

- ▶ NVML\_SUCCESS if version has been set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if vgpuInstance is 0
- ▶ NVML\_ERROR\_NOT\_FOUND if vgpuInstance does not match a valid active vGPU instance on the system
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if length is too small
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Retrieve the NVIDIA driver version installed in the VM associated with a vGPU.

The version is returned as an alphanumeric string in the caller-supplied buffer version. The length of the version string will not exceed 80 characters in length (including the NUL terminator). See [nvmlConstants::NVML\\_SYSTEM\\_DRIVER\\_VERSION\\_BUFFER\\_SIZE](#).

[nvmlVgpuInstanceGetVmDriverVersion\(\)](#) may be called at any time for a vGPU instance. The guest VM driver version is returned as "Not Available" if no NVIDIA driver is installed in the VM, or the VM has not yet booted to the point where the NVIDIA driver is loaded and initialized.

For Kepler or newer fully supported devices.

## **`nvmlReturn_t nvmlVgpuInstanceGetFbUsage (nvmlVgpuInstance_t vgpuInstance, unsigned long long *fbUsage)`**

**Parameters****vgpuInstance**

The identifier of the target instance

**fbUsage**

Pointer to framebuffer usage in bytes

**Returns**

- ▶ NVML\_SUCCESS successful completion

- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `fbUsage` is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the framebuffer usage in bytes.

Framebuffer usage is the amount of vGPU framebuffer memory that is currently in use by the VM.

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuInstanceGetLicenseStatus` (`nvmlVgpuInstance_t vgpuInstance`, unsigned int `*licensed`)

### Parameters

#### `vgpuInstance`

Identifier of the target vGPU instance

#### `licensed`

Reference to return the licensing status

### Returns

- ▶ `NVML_SUCCESS` if `licensed` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `licensed` is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Deprecated Use `nvmlVgpuInstanceGetLicenseInfo_v2`.

Retrieve the current licensing state of the vGPU instance.

If the vGPU is currently licensed, `licensed` is set to 1, otherwise it is set to 0.

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuInstanceGetType (nvmlVgpuInstance_t vgpuInstance, nvmlVgpuTypeId_t *vgpuTypeId)`

### Parameters

#### **vgpuInstance**

Identifier of the target vGPU instance

#### **vgpuTypeId**

Reference to return the vgpuTypeId

### Returns

- ▶ `NVML_SUCCESS` if vgpuTypeId has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if vgpuInstance is 0, or vgpuTypeId is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if vgpuInstance does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the vGPU type of a vGPU instance.

Returns the vGPU type ID of vgpu assigned to the vGPU instance.

For Kepler or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuInstanceGetFrameRateLimit (nvmlVgpuInstance_t vgpuInstance, unsigned int *frameRateLimit)`

### Parameters

#### **vgpuInstance**

Identifier of the target vGPU instance

#### **frameRateLimit**

Reference to return the frame rate limit

### Returns

- ▶ `NVML_SUCCESS` if frameRateLimit has been set

- ▶ `NVML_ERROR_NOT_SUPPORTED` if frame rate limiter is turned off for the vGPU type
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `frameRateLimit` is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the frame rate limit set for the vGPU instance.

Returns the value of the frame rate limit set for the vGPU instance

For Kepler or newer fully supported devices.

```
nvmlReturn_t nvmlVgpuInstanceGetEccMode
(nvmlVgpuInstance_t vgpuInstance, nvmlEnableState_t
*eccMode)
```

### Parameters

#### **vgpuInstance**

The identifier of the target vGPU instance

#### **eccMode**

Reference in which to return the current ECC mode

### Returns

- ▶ `NVML_SUCCESS` if the `vgpuInstance`'s ECC mode has been successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `mode` is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the vGPU doesn't support this feature
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the current ECC mode of vGPU instance.

## `nvmlReturn_t nvmlVgpuInstanceGetEncoderCapacity` (`nvmlVgpuInstance_t vgpuInstance`, unsigned int `*encoderCapacity`)

### Parameters

#### `vgpuInstance`

Identifier of the target vGPU instance

#### `encoderCapacity`

Reference to an unsigned int for the encoder capacity

### Returns

- ▶ `NVML_SUCCESS` if `encoderCapacity` has been retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `encoderQueryType` is invalid
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the encoder capacity of a vGPU instance, as a percentage of maximum encoder capacity with valid values in the range 0-100.

For Maxwell or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuInstanceSetEncoderCapacity` (`nvmlVgpuInstance_t vgpuInstance`, unsigned int `encoderCapacity`)

### Parameters

#### `vgpuInstance`

Identifier of the target vGPU instance

#### `encoderCapacity`

Unsigned int for the encoder capacity value

### Returns

- ▶ `NVML_SUCCESS` if `encoderCapacity` has been set
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized

- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `encoderCapacity` is out of range of 0-100.
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Set the encoder capacity of a vGPU instance, as a percentage of maximum encoder capacity with valid values in the range 0-100.

For Maxwell or newer fully supported devices.

```
nvmlReturn_t nvmlVgpuInstanceGetEncoderStats  
(nvmlVgpuInstance_t vgpuInstance, unsigned int  
*sessionCount, unsigned int *averageFps, unsigned int  
*averageLatency)
```

### Parameters

#### **vgpuInstance**

Identifier of the target vGPU instance

#### **sessionCount**

Reference to an unsigned int for count of active encoder sessions

#### **averageFps**

Reference to an unsigned int for trailing average FPS of all active sessions

#### **averageLatency**

Reference to an unsigned int for encode latency in microseconds

### Returns

- ▶ `NVML_SUCCESS` if `sessionCount`, `averageFps` and `averageLatency` is fetched
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `sessionCount` , or `averageFps` or `averageLatency` is `NULL` or `vgpuInstance` is 0.
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the current encoder statistics of a vGPU Instance

For Maxwell or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuInstanceGetEncoderSessions` (`nvmlVgpuInstance_t vgpuInstance`, `unsigned int *sessionCount`, `nvmlEncoderSessionInfo_t *sessionInfo`)

### Parameters

#### **vgpuInstance**

Identifier of the target vGPU instance

#### **sessionCount**

Reference to caller supplied array size, and returns the number of sessions.

#### **sessionInfo**

Reference to caller supplied array in which the list of session information us returned.

### Returns

- ▶ `NVML_SUCCESS` if `sessionInfo` is fetched
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `sessionCount` is too small, array element count is returned in `sessionCount`
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `sessionCount` is `NULL`, or `vgpuInstance` is 0.
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves information about all active encoder sessions on a vGPU Instance.

An array of active encoder sessions is returned in the caller-supplied buffer pointed at by `sessionInfo`. The array element count is passed in `sessionCount`, and `sessionCount` is used to return the number of sessions written to the buffer.

If the supplied buffer is not large enough to accommodate the active session array, the function returns `NVML_ERROR_INSUFFICIENT_SIZE`, with the element count of `nvmlEncoderSessionInfo_t` array required in `sessionCount`. To query the number of active encoder sessions, call this function with `*sessionCount = 0`. The code will return `NVML_SUCCESS` with number of active encoder sessions updated in `*sessionCount`.

For Maxwell or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuInstanceGetFBCStats` (`nvmlVgpuInstance_t vgpuInstance`, `nvmlFBCStats_t *fbcStats`)

### Parameters

#### `vgpuInstance`

Identifier of the target vGPU instance

#### `fbcStats`

Reference to `nvmlFBCStats_t` structure containing NvFBC stats

### Returns

- ▶ `NVML_SUCCESS` if `fbcStats` is fetched
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `fbcStats` is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the active frame buffer capture sessions statistics of a vGPU Instance

For Maxwell or newer fully supported devices.

## `nvmlReturn_t nvmlVgpuInstanceGetFBCSessions` (`nvmlVgpuInstance_t vgpuInstance`, `unsigned int *sessionCount`, `nvmlFBCSessionInfo_t *sessionInfo`)

### Parameters

#### `vgpuInstance`

Identifier of the target vGPU instance

#### `sessionCount`

Reference to caller supplied array size, and returns the number of sessions.

#### `sessionInfo`

Reference in which to return the session information

### Returns

- ▶ `NVML_SUCCESS` if `sessionInfo` is fetched
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized



- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if vgpuInstance is 0, or sessionCount is NULL.
- ▶ NVML\_ERROR\_NOT\_FOUND if vgpuInstance does not match a valid active vGPU instance on the system
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if sessionCount is too small, array element count is returned in sessionCount
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves information about active frame buffer capture sessions on a vGPU Instance.

An array of active FBC sessions is returned in the caller-supplied buffer pointed at by sessionInfo. The array element count is passed in sessionCount, and sessionCount is used to return the number of sessions written to the buffer.

If the supplied buffer is not large enough to accommodate the active session array, the function returns NVML\_ERROR\_INSUFFICIENT\_SIZE, with the element count of `nvmlFBCSessionInfo_t` array required in sessionCount. To query the number of active FBC sessions, call this function with `*sessionCount = 0`. The code will return NVML\_SUCCESS with number of active FBC sessions updated in `*sessionCount`.

For Maxwell or newer fully supported devices.



hResolution, vResolution, averageFPS and averageLatency data for a FBC session returned in sessionInfo may be zero if there are no new frames captured since the session started.

## nvmlReturn\_t nvmlVgpuInstanceGetGpuInstanceId (nvmlVgpuInstance\_t vgpuInstance, unsigned int \*gpuInstanceId)

### Parameters

#### vgpuInstance

Identifier of the target vGPU instance

#### gpuInstanceId

GPU Instance ID

### Returns

- ▶ NVML\_SUCCESS successful completion
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized

- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `gpuInstanceId` is `NULL`.
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieve the GPU Instance ID for the given vGPU Instance. The API will return a valid GPU Instance ID for MIG backed vGPU Instance, else `INVALID_GPU_INSTANCE_ID` is returned.

For Kepler or newer fully supported devices.

**`nvmlReturn_t nvmlVgpuInstanceGetGpuPciId`**  
**`(nvmlVgpuInstance_t vgpuInstance, char *vgpuPciId,`**  
**`unsigned int *length)`**

### Parameters

#### **`vgpuInstance`**

Identifier of the target vGPU instance

#### **`vgpuPciId`**

Caller-supplied buffer to return vGPU PCI Id string

#### **`length`**

Size of the `vgpuPciId` buffer

### Returns

- ▶ `NVML_SUCCESS` if vGPU PCI Id is successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `vgpuPciId` is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_DRIVER_NOT_LOADED` if NVIDIA driver is not running on the vGPU instance
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `length` is too small, `length` is set to required `length`
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves the PCI Id of the given vGPU Instance i.e. the PCI Id of the GPU as seen inside the VM.

The vGPU PCI id is returned as "00000000:00:00.0" if NVIDIA driver is not installed on the vGPU instance.

## **nvmlReturn\_t nvmlVgpuTypeGetCapabilities (nvmlVgpuTypeId\_t vgpuTypeId, nvmlVgpuCapability\_t capability, unsigned int \*capResult)**

### **Parameters**

#### **vgpuTypeId**

Handle to vGPU type

#### **capability**

Specifies the nvmlVgpuCapability\_t to be queried

#### **capResult**

A boolean for the queried capability indicating that feature is supported

### **Returns**

- ▶ NVML\_SUCCESS successful completion
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if vgpuTypeId is invalid, or capability is invalid, or capResult is NULL
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### **Description**

Retrieve the requested capability for a given vGPU type. Refer to the nvmlVgpuCapability\_t structure for the specific capabilities that can be queried. The return value in capResult should be treated as a boolean, with a non-zero value indicating that the capability is supported.

For Maxwell or newer fully supported devices.

## **nvmlReturn\_t nvmlVgpuInstanceGetMdevUUID (nvmlVgpuInstance\_t vgpuInstance, char \*mdevUuid, unsigned int size)**

### **Parameters**

#### **vgpuInstance**

Identifier of the target vGPU instance

**mdevUuid**

Pointer to caller-supplied buffer to hold MDEV UUID

**size**

Size of buffer in bytes

**Returns**

- ▶ `NVML_SUCCESS` successful completion
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_NOT_SUPPORTED` on any hypervisor other than KVM
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `mdevUuid` is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `size` is too small
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

**Description**

Retrieve the MDEV UUID of a vGPU instance.

The MDEV UUID is a globally unique identifier of the mdev device assigned to the VM, and is returned as a 5-part hexadecimal string, not exceeding 80 characters in length (including the `NULL` terminator). MDEV UUID is displayed only on KVM platform. See [nvmlConstants::NVML\\_DEVICE\\_UUID\\_BUFFER\\_SIZE](#).

For Maxwell or newer fully supported devices.

## 4.25. vGPU Migration

This chapter describes operations that are associated with vGPU Migration.

**struct nvmlVgpuVersion\_t**

**struct nvmlVgpuMetadata\_t**

**struct nvmlVgpuPgpuMetadata\_t**

**struct nvmlVgpuPgpuCompatibility\_t**

**enum nvmlVgpuVmCompatibility\_t**

vGPU VM compatibility codes

#### Values

**NVML\_VGPU\_VM\_COMPATIBILITY\_NONE = 0x0**

vGPU is not runnable

**NVML\_VGPU\_VM\_COMPATIBILITY\_COLD = 0x1**

vGPU is runnable from a cold / powered-off state (ACPI S5)

**NVML\_VGPU\_VM\_COMPATIBILITY\_HIBERNATE = 0x2**

vGPU is runnable from a hibernated state (ACPI S4)

**NVML\_VGPU\_VM\_COMPATIBILITY\_SLEEP = 0x4**

vGPU is runnable from a slept state (ACPI S3)

**NVML\_VGPU\_VM\_COMPATIBILITY\_LIVE = 0x8**

vGPU is runnable from a live/paused (ACPI S0)

**enum nvmlVgpuPgpuCompatibilityLimitCode\_t**

vGPU-pGPU compatibility limit codes

#### Values

**NVML\_VGPU\_COMPATIBILITY\_LIMIT\_NONE = 0x0**

Compatibility is not limited.

**NVML\_VGPU\_COMPATIBILITY\_LIMIT\_HOST\_DRIVER = 0x1**

compatibility is limited by host driver version.

**NVML\_VGPU\_COMPATIBILITY\_LIMIT\_GUEST\_DRIVER = 0x2**

Compatibility is limited by guest driver version.

**NVML\_VGPU\_COMPATIBILITY\_LIMIT\_GPU = 0x4**

Compatibility is limited by GPU hardware.

**NVML\_VGPU\_COMPATIBILITY\_LIMIT\_OTHER = 0x80000000**

Compatibility is limited by an undefined factor.

```

nvmReturn_t nvmVgpuInstanceGetMetadata
(nvmVgpuInstance_t vgpuInstance,
nvmVgpuMetadata_t *vgpuMetadata, unsigned int
*bufferSize)

```

### Parameters

#### **vgpuInstance**

vGPU instance handle

#### **vgpuMetadata**

Pointer to caller-supplied buffer into which vGPU metadata is written

#### **bufferSize**

Size of vgpuMetadata buffer

### Returns

- ▶ NVML\_SUCCESS vGPU metadata structure was successfully returned
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE vgpuMetadata buffer is too small, required size is returned in bufferSize
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if bufferSize is NULL or vgpuInstance is 0; if vgpuMetadata is NULL and the value of bufferSize is not 0.
- ▶ NVML\_ERROR\_NOT\_FOUND if vgpuInstance does not match a valid active vGPU instance on the system
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Returns vGPU metadata structure for a running vGPU. The structure contains information about the vGPU and its associated VM such as the currently installed NVIDIA guest driver version, together with host driver version and an opaque data section containing internal state.

`nvmVgpuInstanceGetMetadata()` may be called at any time for a vGPU instance. Some fields in the returned structure are dependent on information obtained from the guest VM, which may not yet have reached a state where that information is available. The current state of these dependent fields is reflected in the info structure's `nvmVgpuGuestInfoState_t` field.

The VMM may choose to read and save the vGPU's VM info as persistent metadata associated with the VM, and provide it to Virtual GPU Manager when creating a vGPU for subsequent instances of the VM.

The caller passes in a buffer via `vgpuMetadata`, with the size of the buffer in `bufferSize`. If the vGPU Metadata structure is too large to fit in the supplied buffer, the function returns `NVML_ERROR_INSUFFICIENT_SIZE` with the size needed in `bufferSize`.

```
nvmlReturn_t nvmlDeviceGetVgpuMetadata
(nvmlDevice_t device, nvmlVgpuPgpuMetadata_t
*pgpuMetadata, unsigned int *bufferSize)
```

### Parameters

#### **device**

The identifier of the target device

#### **pgpuMetadata**

Pointer to caller-supplied buffer into which `pgpuMetadata` is written

#### **bufferSize**

Pointer to size of `pgpuMetadata` buffer

### Returns

- ▶ `NVML_SUCCESS` GPU metadata structure was successfully returned
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` `pgpuMetadata` buffer is too small, required size is returned in `bufferSize`
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `bufferSize` is `NULL` or device is invalid; if `pgpuMetadata` is `NULL` and the value of `bufferSize` is not 0.
- ▶ `NVML_ERROR_NOT_SUPPORTED` vGPU is not supported by the system
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Returns a vGPU metadata structure for the physical GPU indicated by `device`. The structure contains information about the GPU and the currently installed NVIDIA host driver version that's controlling it, together with an opaque data section containing internal state.

The caller passes in a buffer via `pgpuMetadata`, with the size of the buffer in `bufferSize`. If the `pgpuMetadata` structure is too large to fit in the supplied buffer, the function returns `NVML_ERROR_INSUFFICIENT_SIZE` with the size needed in `bufferSize`.

```
nvmlReturn_t nvmlGetVgpuCompatibility
(nvmlVgpuMetadata_t *vgpuMetadata,
```

## `nvmlVgpuPgpuMetadata_t *pgpuMetadata, nvmlVgpuPgpuCompatibility_t *compatibilityInfo)`

### Parameters

#### **vgpuMetadata**

Pointer to caller-supplied vGPU metadata structure

#### **pgpuMetadata**

Pointer to caller-supplied GPU metadata structure

#### **compatibilityInfo**

Pointer to caller-supplied buffer to hold compatibility info

### Returns

- ▶ NVML\_SUCCESS vGPU metadata structure was successfully returned
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if vgpuMetadata or pgpuMetadata or bufferSize are NULL
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Takes a vGPU instance metadata structure read from `nvmlVgpuInstanceGetMetadata()`, and a vGPU metadata structure for a physical GPU read from `nvmlDeviceGetVgpuMetadata()`, and returns compatibility information of the vGPU instance and the physical GPU.

The caller passes in a buffer via `compatibilityInfo`, into which a compatibility information structure is written. The structure defines the states in which the vGPU / VM may be booted on the physical GPU. If the vGPU / VM compatibility with the physical GPU is limited, a limit code indicates the factor limiting compatibility. (see `nvmlVgpuPgpuCompatibilityLimitCode_t` for details).

Note: vGPU compatibility does not take into account dynamic capacity conditions that may limit a system's ability to boot a given vGPU or associated VM.

## `nvmlReturn_t nvmlDeviceGetPgpuMetadataString (nvmlDevice_t device, char *pgpuMetadata, unsigned int *bufferSize)`

### Parameters

#### **device**

The identifier of the target device



**pgpuMetadata**

Pointer to caller-supplied buffer into which pgpuMetadata is written

**bufferSize**

Pointer to size of pgpuMetadata buffer

**Returns**

- ▶ NVML\_SUCCESS GPU metadata structure was successfully returned
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE pgpuMetadata buffer is too small, required size is returned in bufferSize
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if bufferSize is NULL or device is invalid; if pgpuMetadata is NULL and the value of bufferSize is not 0.
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if vGPU is not supported by the system
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Returns the properties of the physical GPU indicated by the device in an ascii-encoded string format.

The caller passes in a buffer via pgpuMetadata, with the size of the buffer in bufferSize. If the string is too large to fit in the supplied buffer, the function returns NVML\_ERROR\_INSUFFICIENT\_SIZE with the size needed in bufferSize.

## nvmlReturn\_t nvmlDeviceGetVgpuSchedulerLog (nvmlDevice\_t device, nvmlVgpuSchedulerLog\_t \*pSchedulerLog)

**Parameters****device**

The identifier of the target device

**pSchedulerLog**

Reference in which pSchedulerLog is written

**Returns**

- ▶ NVML\_SUCCESS vGPU scheduler logs were successfully obtained
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if pSchedulerLog is NULL or device is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED The API is not supported in current state or device not in vGPU host mode
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Returns the vGPU Software scheduler logs. `pSchedulerLog` points to a caller-allocated structure to contain the logs. The number of elements returned will never exceed `NVML_SCHEDULER_SW_MAX_LOG_ENTRIES`.

To get the entire logs, call the function atleast 5 times a second.

For Pascal or newer fully supported devices.

**`nvmlReturn_t nvmlDeviceGetVgpuSchedulerState`  
(`nvmlDevice_t device`, `nvmlVgpuSchedulerGetState_t *pSchedulerState`)**

**Parameters****device**

The identifier of the target device

**pSchedulerState**

Reference in which `pSchedulerState` is returned

**Returns**

- ▶ `NVML_SUCCESS` vGPU scheduler state is successfully obtained
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `pSchedulerState` is NULL or device is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` The API is not supported in current state or device not in vGPU host mode
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

**Description**

Returns the vGPU scheduler state. The information returned in `nvmlVgpuSchedulerGetState_t` is not relevant if the BEST EFFORT policy is set.

For Pascal or newer fully supported devices.

## nvmlReturn\_t nvmlDeviceGetVgpuSchedulerCapabilities (nvmlDevice\_t device, nvmlVgpuSchedulerCapabilities\_t \*pCapabilities)

### Parameters

#### device

The identifier of the target device

#### pCapabilities

Reference in which pCapabilities is written

### Returns

- ▶ NVML\_SUCCESS vGPU scheduler capabilities were successfully obtained
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if pCapabilities is NULL or device is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED The API is not supported in current state or device not in vGPU host mode
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Returns the vGPU scheduler capabilities. The list of supported vGPU schedulers returned in `nvmlVgpuSchedulerCapabilities_t` is from the `NVML_VGPU_SCHEDULER_POLICY_*`. This list enumerates the supported scheduler policies if the engine is Graphics type. The other values in `nvmlVgpuSchedulerCapabilities_t` are also applicable if the engine is Graphics type. For other engine types, it is BEST EFFORT policy. If ARR is supported and enabled, scheduling frequency and averaging factor are applicable else timeSlice is applicable.

For Pascal or newer fully supported devices.

## nvmlReturn\_t nvmlDeviceSetVgpuSchedulerState (nvmlDevice\_t device, nvmlVgpuSchedulerSetState\_t \*pSchedulerState)

### Parameters

#### device

The identifier of the target device

#### pSchedulerState

vGPU pSchedulerState to set

**Returns**

- ▶ NVML\_SUCCESS vGPU scheduler state has been successfully set
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if pSchedulerState is NULL or device is invalid
- ▶ NVML\_ERROR\_RESET\_REQUIRED if setting pSchedulerState failed with fatal error, reboot is required to overcome from this error.
- ▶ NVML\_ERROR\_NOT\_SUPPORTED The API is not supported in current state or device not in vGPU host mode or if any vGPU instance currently exists on the device
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Sets the vGPU scheduler state.

For Pascal or newer fully supported devices.

The scheduler state change won't persist across module load/unload. Scheduler state and params will be allowed to set only when no VM is running. In `nvmlVgpuSchedulerSetState_t`, IFF `enableARRMode` is enabled then provide `avgFactorForARR` and frequency as input. If `enableARRMode` is disabled then provide `timeslice` as input.

## `nvmlReturn_t nvmlGetVgpuVersion (nvmlVgpuVersion_t *supported, nvmlVgpuVersion_t *current)`

**Parameters****supported**

Pointer to the structure in which the preset range of vGPU versions supported by the NVIDIA vGPU Manager is written

**current**

Pointer to the structure in which the range of supported vGPU versions set by an administrator is written

**Returns**

- ▶ NVML\_SUCCESS The vGPU version range structures were successfully obtained.
- ▶ NVML\_ERROR\_NOT\_SUPPORTED The API is not supported.
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT The supported parameter or the current parameter is NULL.
- ▶ NVML\_ERROR\_UNKNOWN An error occurred while the data was being fetched.

## Description

Query the ranges of supported vGPU versions.

This function gets the linear range of supported vGPU versions that is preset for the NVIDIA vGPU Manager and the range set by an administrator. If the preset range has not been overridden by `nvmlSetVgpuVersion`, both ranges are the same.

The caller passes pointers to the following `nvmlVgpuVersion_t` structures, into which the NVIDIA vGPU Manager writes the ranges: 1. supported structure that represents the preset range of vGPU versions supported by the NVIDIA vGPU Manager. 2. current structure that represents the range of supported vGPU versions set by an administrator. By default, this range is the same as the preset range.

## `nvmlReturn_t nvmlSetVgpuVersion (nvmlVgpuVersion_t *vgpuVersion)`

### Parameters

#### `vgpuVersion`

Pointer to a caller-supplied range of supported vGPU versions.

### Returns

- ▶ `NVML_SUCCESS` The preset range of supported vGPU versions was successfully overridden.
- ▶ `NVML_ERROR_NOT_SUPPORTED` The API is not supported.
- ▶ `NVML_ERROR_IN_USE` The range was not overridden because a VM is running on the host.
- ▶ `NVML_ERROR_INVALID_ARGUMENT` The `vgpuVersion` parameter specifies a range that is outside the range supported by the NVIDIA vGPU Manager or if `vgpuVersion` is `NULL`.

## Description

Override the preset range of vGPU versions supported by the NVIDIA vGPU Manager with a range set by an administrator.

This function configures the NVIDIA vGPU Manager with a range of supported vGPU versions set by an administrator. This range must be a subset of the preset range that the NVIDIA vGPU Manager supports. The custom range set by an administrator takes precedence over the preset range and is advertised to the guest VM for negotiating the vGPU version. See `nvmlGetVgpuVersion` for details of how to query the preset range of versions supported.

This function takes a pointer to vGPU version range structure `nvmlVgpuVersion_t` as input to override the preset vGPU version range that the NVIDIA vGPU Manager supports.

After host system reboot or driver reload, the range of supported versions reverts to the range that is preset for the NVIDIA vGPU Manager.



1. The range set by the administrator must be a subset of the preset range that the NVIDIA vGPU Manager supports. Otherwise, an error is returned. 2. If the range of supported guest driver versions does not overlap the range set by the administrator, the guest driver fails to load. 3. If the range of supported guest driver versions overlaps the range set by the administrator, the guest driver will load with a negotiated vGPU version that is the maximum value in the overlapping range. 4. No VMs must be running on the host when this function is called. If a VM is running on the host, the call to this function fails.

## 4.26. vGPU Utilization and Accounting

This chapter describes operations that are associated with vGPU Utilization and Accounting.

```
nvmlReturn_t nvmlDeviceGetVgpuUtilization
(nvmlDevice_t device, unsigned long long
lastSeenTimeStamp, nvmlValueType_t *sampleValType,
unsigned int *vgpuInstanceSamplesCount,
nvmlVgpuInstanceUtilizationSample_t
*utilizationSamples)
```

### Parameters

#### **device**

The identifier for the target device

#### **lastSeenTimeStamp**

Return only samples with timestamp greater than lastSeenTimeStamp.

#### **sampleValType**

Pointer to caller-supplied buffer to hold the type of returned sample values

#### **vgpuInstanceSamplesCount**

Pointer to caller-supplied array size, and returns number of vGPU instances

#### **utilizationSamples**

Pointer to caller-supplied buffer in which vGPU utilization samples are returned

## Returns

- ▶ NVML\_SUCCESS if utilization samples are successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, vgpuInstanceSamplesCount or sampleValType is NULL, or a sample count of 0 is passed with a non-NULL utilizationSamples
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if supplied vgpuInstanceSamplesCount is too small to return samples for all vGPU instances currently executing on the device
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if vGPU is not supported by the device
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_NOT\_FOUND if sample entries are not found
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

## Description

Retrieves current utilization for vGPUs on a physical GPU (device).

For Kepler or newer fully supported devices.

Reads recent utilization of GPU SM (3D/Compute), framebuffer, video encoder, and video decoder for vGPU instances running on a device. Utilization values are returned as an array of utilization sample structures in the caller-supplied buffer pointed at by utilizationSamples. One utilization sample structure is returned per vGPU instance, and includes the CPU timestamp at which the samples were recorded. Individual utilization values are returned as "unsigned int" values in `nvmlValue_t` unions. The function sets the caller-supplied sampleValType to NVML\_VALUE\_TYPE\_UNSIGNED\_INT to indicate the returned value type.

To read utilization values, first determine the size of buffer required to hold the samples by invoking the function with utilizationSamples set to NULL. The function will return NVML\_ERROR\_INSUFFICIENT\_SIZE, with the current vGPU instance count in vgpuInstanceSamplesCount, or NVML\_SUCCESS if the current vGPU instance count is zero. The caller should allocate a buffer of size vgpuInstanceSamplesCount \* sizeof(nvmlVgpuInstanceUtilizationSample\_t). Invoke the function again with the allocated buffer passed in utilizationSamples, and vgpuInstanceSamplesCount set to the number of entries the buffer is sized for.

On successful return, the function updates vgpuInstanceSampleCount with the number of vGPU utilization sample structures that were actually written. This may differ from a previously read value as vGPU instances are created or destroyed.

lastSeenTimeStamp represents the CPU timestamp in microseconds at which utilization samples were last read. Set it to 0 to read utilization based on all the samples maintained by the driver's internal sample buffer. Set lastSeenTimeStamp to a timeStamp retrieved from a previous query to read utilization since the previous query.

## nvmlReturn\_t nvmlDeviceGetVgpuInstancesUtilizationInfo (nvmlDevice\_t device, nvmlVgpuInstancesUtilizationInfo\_t \*vgpuUtilInfo)

### Parameters

#### device

The identifier for the target device

#### vgpuUtilInfo

Pointer to the caller-provided structure of nvmlVgpuInstancesUtilizationInfo\_t

### Returns

- ▶ NVML\_SUCCESS if utilization samples are successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, vgpuUtilInfo is NULL, or vgpuUtilInfo->vgpuInstanceCount is 0
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if vGPU is not supported by the device
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_VERSION\_MISMATCH if the version of vgpuUtilInfo is invalid
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if vgpuUtilInfo->vgpuUtilArray is NULL, or the buffer size of vgpuUtilInfo->vgpuInstanceCount is too small. The caller should check the current vGPU instance count from the returned vgpuUtilInfo->vgpuInstanceCount, and call the function again with a buffer of size vgpuUtilInfo->vgpuInstanceCount \* sizeof(nvmlVgpuInstanceUtilizationInfo\_t)
- ▶ NVML\_ERROR\_NOT\_FOUND if sample entries are not found
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Retrieves recent utilization for vGPU instances running on a physical GPU (device).

For Kepler or newer fully supported devices.

Reads recent utilization of GPU SM (3D/Compute), framebuffer, video encoder, video decoder, jpeg decoder, and OFA for vGPU instances running on a device. Utilization values are returned as an array of utilization sample structures in the caller-supplied buffer pointed at by vgpuUtilInfo->vgpuUtilArray. One utilization sample structure is returned per vGPU instance, and includes the CPU timestamp at which the samples were recorded. Individual utilization values are returned as "unsigned int" values in



`nvmlValue_t` unions. The function sets the caller-supplied `vgpuUtilInfo->sampleValType` to `NVML_VALUE_TYPE_UNSIGNED_INT` to indicate the returned value type.

To read utilization values, first determine the size of buffer required to hold the samples by invoking the function with `vgpuUtilInfo->vgpuUtilArray` set to `NULL`. The function will return `NVML_ERROR_INSUFFICIENT_SIZE`, with the current vGPU instance count in `vgpuUtilInfo->vgpuInstanceCount`, or `NVML_SUCCESS` if the current vGPU instance count is zero. The caller should allocate a buffer of size `vgpuUtilInfo->vgpuInstanceCount * sizeof(nvmlVgpuInstanceUtilizationInfo_t)`. Invoke the function again with the allocated buffer passed in `vgpuUtilInfo->vgpuUtilArray`, and `vgpuUtilInfo->vgpuInstanceCount` set to the number of entries the buffer is sized for.

On successful return, the function updates `vgpuUtilInfo->vgpuInstanceCount` with the number of vGPU utilization sample structures that were actually written. This may differ from a previously read value as vGPU instances are created or destroyed.

`vgpuUtilInfo->lastSeenTimeStamp` represents the CPU timestamp in microseconds at which utilization samples were last read. Set it to 0 to read utilization based on all the samples maintained by the driver's internal sample buffer. Set `vgpuUtilInfo->lastSeenTimeStamp` to a `timeStamp` retrieved from a previous query to read utilization since the previous query.

**`nvmlReturn_t nvmlDeviceGetVgpuProcessUtilization`**  
**(`nvmlDevice_t device`, unsigned long**  
**long `lastSeenTimeStamp`, unsigned**  
**int `*vgpuProcessSamplesCount`,**  
**`nvmlVgpuProcessUtilizationSample_t`**  
**`*utilizationSamples`)**

#### Parameters

##### **device**

The identifier for the target device

##### **lastSeenTimeStamp**

Return only samples with timestamp greater than `lastSeenTimeStamp`.

##### **vgpuProcessSamplesCount**

Pointer to caller-supplied array size, and returns number of processes running on vGPU instances

##### **utilizationSamples**

Pointer to caller-supplied buffer in which vGPU sub process utilization samples are returned

## Returns

- ▶ NVML\_SUCCESS if utilization samples are successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid, vgpuProcessSamplesCount or a sample count of 0 is passed with a non-NULL utilizationSamples
- ▶ NVML\_ERROR\_INSUFFICIENT\_SIZE if supplied vgpuProcessSamplesCount is too small to return samples for all vGPU instances currently executing on the device
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if vGPU is not supported by the device
- ▶ NVML\_ERROR\_GPU\_IS\_LOST if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ NVML\_ERROR\_NOT\_FOUND if sample entries are not found
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

## Description

Retrieves current utilization for processes running on vGPUs on a physical GPU (device).

For Maxwell or newer fully supported devices.

Reads recent utilization of GPU SM (3D/Compute), framebuffer, video encoder, and video decoder for processes running on vGPU instances active on a device. Utilization values are returned as an array of utilization sample structures in the caller-supplied buffer pointed at by utilizationSamples. One utilization sample structure is returned per process running on vGPU instances, that had some non-zero utilization during the last sample period. It includes the CPU timestamp at which the samples were recorded. Individual utilization values are returned as "unsigned int" values.

To read utilization values, first determine the size of buffer required to hold the samples by invoking the function with utilizationSamples set to NULL. The function will return NVML\_ERROR\_INSUFFICIENT\_SIZE, with the current vGPU instance count in vgpuProcessSamplesCount. The caller should allocate a buffer of size vgpuProcessSamplesCount \* sizeof(nvmlVgpuProcessUtilizationSample\_t). Invoke the function again with the allocated buffer passed in utilizationSamples, and vgpuProcessSamplesCount set to the number of entries the buffer is sized for.

On successful return, the function updates vgpuSubProcessSampleCount with the number of vGPU sub process utilization sample structures that were actually written. This may differ from a previously read value depending on the number of processes that are active in any given sample period.

lastSeenTimeStamp represents the CPU timestamp in microseconds at which utilization samples were last read. Set it to 0 to read utilization based on all the samples maintained by the driver's internal sample buffer. Set lastSeenTimeStamp to a timeStamp retrieved from a previous query to read utilization since the previous query.

## `nvmlReturn_t` `nvmlDeviceGetVgpuProcessesUtilizationInfo` (`nvmlDevice_t` device, `nvmlVgpuProcessesUtilizationInfo_t *vgpuProcUtilInfo`)

### Parameters

#### **device**

The identifier for the target device

#### **vgpuProcUtilInfo**

Pointer to the caller-provided structure of `nvmlVgpuProcessesUtilizationInfo_t`

### Returns

- ▶ `NVML_SUCCESS` if utilization samples are successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device is invalid, or `vgpuProcUtilInfo` is null
- ▶ `NVML_ERROR_VERSION_MISMATCH` if the version of `vgpuProcUtilInfo` is invalid
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `vgpuProcUtilInfo->vgpuProcUtilArray` is null, or supplied `vgpuProcUtilInfo->vgpuProcessCount` is too small to return samples for all processes on vGPU instances currently executing on the device. The caller should check the current processes count from the returned `vgpuProcUtilInfo->vgpuProcessCount`, and call the function again with a buffer of size `vgpuProcUtilInfo->vgpuProcessCount * sizeof(nvmlVgpuProcessUtilizationSample_t)`
- ▶ `NVML_ERROR_NOT_SUPPORTED` if vGPU is not supported by the device
- ▶ `NVML_ERROR_GPU_IS_LOST` if the target GPU has fallen off the bus or is otherwise inaccessible
- ▶ `NVML_ERROR_NOT_FOUND` if sample entries are not found
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Retrieves recent utilization for processes running on vGPU instances on a physical GPU (device).

For Maxwell or newer fully supported devices.

Reads recent utilization of GPU SM (3D/Compute), framebuffer, video encoder, video decoder, jpeg decoder, and OFA for processes running on vGPU instances active on a device. Utilization values are returned as an array of utilization sample structures in

the caller-supplied buffer pointed at by `vgpuProcUtilInfo->vgpuProcUtilArray`. One utilization sample structure is returned per process running on vGPU instances, that had some non-zero utilization during the last sample period. It includes the CPU timestamp at which the samples were recorded. Individual utilization values are returned as "unsigned int" values.

To read utilization values, first determine the size of buffer required to hold the samples by invoking the function with `vgpuProcUtilInfo->vgpuProcUtilArray` set to NULL. The function will return `NVML_ERROR_INSUFFICIENT_SIZE`, with the current processes' count running on vGPU instances in `vgpuProcUtilInfo->vgpuProcessCount`. The caller should allocate a buffer of size `vgpuProcUtilInfo->vgpuProcessCount * sizeof(nvmlVgpuProcessUtilizationSample_t)`. Invoke the function again with the allocated buffer passed in `vgpuProcUtilInfo->vgpuProcUtilArray`, and `vgpuProcUtilInfo->vgpuProcessCount` set to the number of entries the buffer is sized for.

On successful return, the function updates `vgpuProcUtilInfo->vgpuProcessCount` with the number of vGPU sub process utilization sample structures that were actually written. This may differ from a previously read value depending on the number of processes that are active in any given sample period.

`vgpuProcUtilInfo->lastSeenTimeStamp` represents the CPU timestamp in microseconds at which utilization samples were last read. Set it to 0 to read utilization based on all the samples maintained by the driver's internal sample buffer. Set `vgpuProcUtilInfo->lastSeenTimeStamp` to a `timeStamp` retrieved from a previous query to read utilization since the previous query.

## `nvmlReturn_t nvmlVgpuInstanceGetAccountingMode(nvmlVgpuInstance_t vgpuInstance, nvmlEnableState_t *mode)`

### Parameters

#### **vgpuInstance**

The identifier of the target vGPU instance

#### **mode**

Reference in which to return the current accounting mode

### Returns

- ▶ `NVML_SUCCESS` if the mode has been successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `mode` is NULL
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system

- ▶ `NVML_ERROR_NOT_SUPPORTED` if the vGPU doesn't support this feature
- ▶ `NVML_ERROR_DRIVER_NOT_LOADED` if NVIDIA driver is not running on the vGPU instance
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Queries the state of per process accounting mode on vGPU.

For Maxwell or newer fully supported devices.

`nvmlReturn_t nvmlVgpuInstanceGetAccountingPids  
(nvmlVgpuInstance_t vgpuInstance, unsigned int *count,  
unsigned int *pids)`

### Parameters

#### **vgpuInstance**

The identifier of the target vGPU instance

#### **count**

Reference in which to provide the pids array size, and to return the number of elements ready to be queried

#### **pids**

Reference in which to return list of process ids

### Returns

- ▶ `NVML_SUCCESS` if pids were successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `count` is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the vGPU doesn't support this feature or accounting mode is disabled
- ▶ `NVML_ERROR_INSUFFICIENT_SIZE` if `count` is too small (`count` is set to expected value)
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Queries list of processes running on vGPU that can be queried for accounting stats. The list of processes returned can be in running or terminated state.

For Maxwell or newer fully supported devices.

To just query the maximum number of processes that can be queried, call this function with `*count = 0` and `pids=NULL`. The return code will be `NVML_ERROR_INSUFFICIENT_SIZE`, or `NVML_SUCCESS` if list is empty.

For more details see [nvmlVgpuInstanceGetAccountingStats](#).



In case of PID collision some processes might not be accessible before the circular buffer is full.

#### See also:

[nvmlVgpuInstanceGetAccountingPids](#)

## `nvmlReturn_t nvmlVgpuInstanceGetAccountingStats` (`nvmlVgpuInstance_t vgpuInstance`, unsigned int `pid`, `nvmlAccountingStats_t *stats`)

### Parameters

#### **vgpuInstance**

The identifier of the target vGPU instance

#### **pid**

Process Id of the target process to query stats for

#### **stats**

Reference in which to return the process's accounting stats

### Returns

- ▶ `NVML_SUCCESS` if stats have been successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is 0, or `stats` is `NULL`
- ▶ `NVML_ERROR_NOT_FOUND` if `vgpuInstance` does not match a valid active vGPU instance on the system or `stats` is not found
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the vGPU doesn't support this feature or accounting mode is disabled
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Queries process's accounting stats.

For Maxwell or newer fully supported devices.

Accounting stats capture GPU utilization and other statistics across the lifetime of a process, and can be queried during life time of the process or after its termination. The time field in `nvmlAccountingStats_t` is reported as 0 during the lifetime of the process and updated to actual running time after its termination. Accounting stats are kept in a circular buffer, newly created processes overwrite information about old processes.

See `nvmlAccountingStats_t` for description of each returned metric. List of processes that can be queried can be retrieved from `nvmlVgpuInstanceGetAccountingPids`.



- ▶ Accounting Mode needs to be on. See `nvmlVgpuInstanceGetAccountingMode`.
- ▶ Only compute and graphics applications stats can be queried. Monitoring applications stats can't be queried since they don't contribute to GPU utilization.
- ▶ In case of pid collision stats of only the latest process (that terminated last) will be reported

## `nvmlReturn_t nvmlVgpuInstanceClearAccountingPids` (`nvmlVgpuInstance_t vgpuInstance`)

### Parameters

#### `vgpuInstance`

The identifier of the target vGPU instance

### Returns

- ▶ `NVML_SUCCESS` if accounting information has been cleared
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `vgpuInstance` is invalid
- ▶ `NVML_ERROR_NO_PERMISSION` if the user doesn't have permission to perform this operation
- ▶ `NVML_ERROR_NOT_SUPPORTED` if the vGPU doesn't support this feature or accounting mode is disabled
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Clears accounting information of the vGPU instance that have already terminated.

For Maxwell or newer fully supported devices. Requires root/admin permissions.



- ▶ Accounting Mode needs to be on. See `nvmlVgpuInstanceGetAccountingMode`.
- ▶ Only compute and graphics applications stats are reported and can be cleared since monitoring applications stats don't contribute to GPU utilization.

```

nvmlReturn_t nvmlVgpuInstanceGetLicenseInfo_v2
(nvmlVgpuInstance_t vgpuInstance,
nvmlVgpuLicenseInfo_t *licenseInfo)

```

### Parameters

#### **vgpuInstance**

Identifier of the target vGPU instance

#### **licenseInfo**

Pointer to vGPU license information structure

### Returns

- ▶ NVML\_SUCCESS if information is successfully retrieved
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if vgpuInstance is 0, or licenseInfo is NULL
- ▶ NVML\_ERROR\_NOT\_FOUND if vgpuInstance does not match a valid active vGPU instance on the system
- ▶ NVML\_ERROR\_DRIVER\_NOT\_LOADED if NVIDIA driver is not running on the vGPU instance
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

### Description

Query the license information of the vGPU instance.

For Maxwell or newer fully supported devices.

## 4.27. Excluded GPU Queries

This chapter describes NVML operations that are associated with excluded GPUs.

```

struct nvmlExcludedDeviceInfo_t

```

```

nvmlReturn_t nvmlGetExcludedDeviceCount (unsigned
int *deviceCount)

```

### Parameters

#### **deviceCount**

Reference in which to return the number of excluded devices



**Returns**

- ▶ NVML\_SUCCESS if deviceCount has been set
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if deviceCount is NULL

**Description**

Retrieves the number of excluded GPU devices in the system.

For all products.

## `nvmlReturn_t nvmlGetExcludedDeviceInfoByIndex` (unsigned int index, nvmlExcludedDeviceInfo\_t \*info)

**Parameters****index**

The index of the target GPU,  $\geq 0$  and  $<$  deviceCount

**info**

Reference in which to return the device information

**Returns**

- ▶ NVML\_SUCCESS if device has been set
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if index is invalid or info is NULL

**Description**

Acquire the device information for an excluded GPU device, based on its index.

For all products.

Valid indices are derived from the deviceCount returned by `nvmlGetExcludedDeviceCount()`. For example, if deviceCount is 2 the valid indices are 0 and 1, corresponding to GPU 0 and GPU 1.

**See also:**

`nvmlGetExcludedDeviceCount`

## 4.28. Multi Instance GPU Management

This chapter describes NVML operations that are associated with Multi Instance GPU management.

```

struct nvmlGpuInstancePlacement_t
struct nvmlGpuInstanceProfileInfo_t
struct nvmlGpuInstanceProfileInfo_v2_t
struct nvmlGpuInstanceProfileInfo_v3_t
struct nvmlComputeInstanceProfileInfo_t
struct nvmlComputeInstanceProfileInfo_v2_t
struct nvmlComputeInstanceProfileInfo_v3_t
nvmlReturn_t nvmlDeviceSetMigMode (nvmlDevice_t
device, unsigned int mode, nvmlReturn_t
*activationStatus)

```

#### Parameters

##### device

The identifier of the target device

##### mode

The mode to be set, `NVML_DEVICE_MIG_DISABLE` or `NVML_DEVICE_MIG_ENABLE`

##### activationStatus

The activationStatus status

#### Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device, mode or activationStatus are invalid
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't support MIG mode

## Description

Set MIG mode for the device.

For Ampere or newer fully supported devices. Requires root user.

This mode determines whether a GPU instance can be created.

This API may unbind or reset the device to activate the requested mode. Thus, the attributes associated with the device, such as minor number, might change. The caller of this API is expected to query such attributes again.

On certain platforms like pass-through virtualization, where reset functionality may not be exposed directly, VM reboot is required. `activationStatus` would return `NVML_ERROR_RESET_REQUIRED` for such cases.

`activationStatus` would return the appropriate error code upon unsuccessful activation. For example, if device unbind fails because the device isn't idle, `NVML_ERROR_IN_USE` would be returned. The caller of this API is expected to idle the device and retry setting the mode.



On Windows, only disabling MIG mode is supported. `activationStatus` would return `NVML_ERROR_NOT_SUPPORTED` as GPU reset is not supported on Windows through this API.

`nvmlReturn_t nvmlDeviceGetMigMode (nvmlDevice_t device, unsigned int *currentMode, unsigned int *pendingMode)`

## Parameters

### **device**

The identifier of the target device

### **currentMode**

Returns the current mode, `NVML_DEVICE_MIG_DISABLE` or `NVML_DEVICE_MIG_ENABLE`

### **pendingMode**

Returns the pending mode, `NVML_DEVICE_MIG_DISABLE` or `NVML_DEVICE_MIG_ENABLE`

## Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized

- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device, `currentMode` or `pendingMode` are invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't support MIG mode

### Description

Get MIG mode for the device.

For Ampere or newer fully supported devices.

Changing MIG modes may require device unbind or reset. The "pending" MIG mode refers to the target mode following the next activation trigger.

```
nvmlReturn_t nvmlDeviceGetGpuInstanceProfileInfo  
(nvmlDevice_t device, unsigned int profile,  
nvmlGpuInstanceProfileInfo_t *info)
```

### Parameters

#### **device**

The identifier of the target device

#### **profile**

One of the `NVML_GPU_INSTANCE_PROFILE_*`

#### **info**

Returns detailed profile information

### Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device, profile or info are invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't support MIG or profile isn't supported
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation

### Description

Get GPU instance profile information

Information provided by this API is immutable throughout the lifetime of a MIG mode.

For Ampere or newer fully supported devices. Supported on Linux only.

## `nvmlReturn_t nvmlDeviceGetGpuInstanceProfileInfoV` (`nvmlDevice_t device`, `unsigned int profile`, `nvmlGpuInstanceProfileInfo_v2_t *info`)

### Parameters

#### **device**

The identifier of the target device

#### **profile**

One of the `NVML_GPU_INSTANCE_PROFILE_*`

#### **info**

Returns detailed profile information

### Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device, profile, info, or info->version are invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't have MIG mode enabled or profile isn't supported
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation

### Description

Versioned wrapper around `nvmlDeviceGetGpuInstanceProfileInfo` that accepts a versioned `nvmlGpuInstanceProfileInfo_v2_t` or later output structure.



The caller must set the `nvmlGpuInstanceProfileInfo_v2_t::version` field to the appropriate version prior to calling this function. For example:

```
nvmlGpuInstanceProfileInfo_v2_t profileInfo =
    { .version = nvmlGpuInstanceProfileInfo_v2 };
nvmlReturn_t result
= nvmlDeviceGetGpuInstanceProfileInfoV(device,
profile,
&profileInfo);
```

For Ampere or newer fully supported devices. Supported on Linux only.

## `nvmlReturn_t` `nvmlDeviceGetGpuInstancePossiblePlacements_v2`

(`nvmlDevice_t device`, unsigned int `profileId`,  
`nvmlGpuInstancePlacement_t *placements`, unsigned int  
`*count`)

### Parameters

#### **device**

The identifier of the target device

#### **profileId**

The GPU instance profile ID. See [nvmlDeviceGetGpuInstanceProfileInfo](#)

#### **placements**

Returns placements allowed for the profile. Can be NULL to discover number of allowed placements for this profile. If non-NULL must be large enough to accommodate the placements supported by the profile.

#### **count**

Returns number of allowed placemenets for the profile.

### Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device, profileId or count are invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't support MIG or profileId isn't supported
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation

### Description

Get GPU instance placements.

A placement represents the location of a GPU instance within a device. This API only returns all the possible placements for the given profile regardless of whether MIG is enabled or not. A created GPU instance occupies memory slices described by its placement. Creation of new GPU instance will fail if there is overlap with the already occupied memory slices.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

`nvmlReturn_t`  
`nvmlDeviceGetGpuInstanceRemainingCapacity`

**(nvmlDevice\_t device, unsigned int profileId, unsigned int \*count)**

### Parameters

#### device

The identifier of the target device

#### profileId

The GPU instance profile ID. See [nvmlDeviceGetGpuInstanceProfileInfo](#)

#### count

Returns remaining instance count for the profile ID

### Returns

- ▶ NVML\_SUCCESS Upon success
- ▶ NVML\_ERROR\_UNINITIALIZED If library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT If device, profileId or count are invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED If device doesn't have MIG mode enabled or profileId isn't supported
- ▶ NVML\_ERROR\_NO\_PERMISSION If user doesn't have permission to perform the operation

### Description

Get GPU instance profile capacity.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

**nvmlReturn\_t nvmlDeviceCreateGpuInstance  
(nvmlDevice\_t device, unsigned int profileId,  
nvmlGpuInstance\_t \*gpuInstance)**

### Parameters

#### device

The identifier of the target device

#### profileId

The GPU instance profile ID. See [nvmlDeviceGetGpuInstanceProfileInfo](#)

#### gpuInstance

Returns the GPU instance handle

**Returns**

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device, profile, profileId or gpuInstance are invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't have MIG mode enabled or in vGPU guest
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation
- ▶ `NVML_ERROR_INSUFFICIENT_RESOURCES` If the requested GPU instance could not be created

**Description**

Create GPU instance.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

If the parent device is unbound, reset or the GPU instance is destroyed explicitly, the GPU instance handle would become invalid. The GPU instance must be recreated to acquire a valid handle.

```
nvmlReturn_t  
nvmlDeviceCreateGpuInstanceWithPlacement  
(nvmlDevice_t device, unsigned int profileId,  
const nvmlGpuInstancePlacement_t *placement,  
nvmlGpuInstance_t *gpuInstance)
```

**Parameters****device**

The identifier of the target device

**profileId**

The GPU instance profile ID. See [nvmlDeviceGetGpuInstanceProfileInfo](#)

**placement**

The requested placement. See [nvmlDeviceGetGpuInstancePossiblePlacements\\_v2](#)

**gpuInstance**

Returns the GPU instance handle

**Returns**

- ▶ `NVML_SUCCESS` Upon success



- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If device, profile, profileId, placement or `gpuInstance` are invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't have MIG mode enabled or in vGPU guest
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation
- ▶ `NVML_ERROR_INSUFFICIENT_RESOURCES` If the requested GPU instance could not be created

### Description

Create GPU instance with the specified placement.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

If the parent device is unbound, reset or the GPU instance is destroyed explicitly, the GPU instance handle would become invalid. The GPU instance must be recreated to acquire a valid handle.

## `nvmlReturn_t nvmlGpuInstanceDestroy` (`nvmlGpuInstance_t gpuInstance`)

### Parameters

#### `gpuInstance`

The GPU instance handle

### Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If `gpuInstance` is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't have MIG mode enabled or in vGPU guest
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation
- ▶ `NVML_ERROR_IN_USE` If the GPU instance is in use. This error would be returned if processes (e.g. CUDA application) or compute instances are active on the GPU instance.

### Description

Destroy GPU instance.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

**nvmlReturn\_t nvmlDeviceGetGpuInstances**  
 (nvmlDevice\_t device, unsigned int profileId,  
 nvmlGpuInstance\_t \*gpuInstances, unsigned int \*count)

### Parameters

#### device

The identifier of the target device

#### profileId

The GPU instance profile ID. See [nvmlDeviceGetGpuInstanceProfileInfo](#)

#### gpuInstances

Returns pre-existing GPU instances, the buffer must be large enough to accommodate the instances supported by the profile. See [nvmlDeviceGetGpuInstanceProfileInfo](#)

#### count

The count of returned GPU instances

### Returns

- ▶ NVML\_SUCCESS Upon success
- ▶ NVML\_ERROR\_UNINITIALIZED If library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT If device, profileId, gpuInstances or count are invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED If device doesn't have MIG mode enabled
- ▶ NVML\_ERROR\_NO\_PERMISSION If user doesn't have permission to perform the operation

### Description

Get GPU instances for given profile ID.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

## `nvmlReturn_t nvmlDeviceGetGpuInstanceById` (`nvmlDevice_t device`, `unsigned int id`, `nvmlGpuInstance_t *gpuInstance`)

### Parameters

#### **device**

The identifier of the target device

#### **id**

The GPU instance ID

#### **gpuInstance**

Returns GPU instance

### Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If `device`, `id` or `gpuInstance` are invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't have MIG mode enabled
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation
- ▶ `NVML_ERROR_NOT_FOUND` If the GPU instance is not found.

### Description

Get GPU instances for given instance ID.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

## `nvmlReturn_t nvmlGpuInstanceGetInfo` (`nvmlGpuInstance_t gpuInstance`, `nvmlGpuInstanceInfo_t *info`)

### Parameters

#### **gpuInstance**

The GPU instance handle

#### **info**

Return GPU instance information

**Returns**

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If `gpuInstance` or `info` are invalid
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation

**Description**

Get GPU instance information.

For Ampere or newer fully supported devices. Supported on Linux only.

```
nvmlReturn_t
nvmlGpuInstanceGetComputeInstanceProfileInfo
(nvmlGpuInstance_t gpuInstance, unsigned
int profile, unsigned int engProfile,
nvmlComputeInstanceProfileInfo_t *info)
```

**Parameters****gpuInstance**

The identifier of the target GPU instance

**profile**

One of the `NVML_COMPUTE_INSTANCE_PROFILE_*`

**engProfile**

One of the `NVML_COMPUTE_INSTANCE_ENGINE_PROFILE_*`

**info**

Returns detailed profile information

**Returns**

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If `gpuInstance`, `profile`, `engProfile` or `info` are invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` If `profile` isn't supported
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation

**Description**

Get compute instance profile information.

Information provided by this API is immutable throughout the lifetime of a MIG mode.  
For Ampere or newer fully supported devices. Supported on Linux only.

**nvmlReturn\_t**  
**nvmlGpuInstanceGetComputeInstanceProfileInfoV**  
 (nvmlGpuInstance\_t gpuInstance, unsigned  
 int profile, unsigned int engProfile,  
 nvmlComputeInstanceProfileInfo\_v2\_t \*info)

### Parameters

#### gpuInstance

The identifier of the target GPU instance

#### profile

One of the NVML\_COMPUTE\_INSTANCE\_PROFILE\_\*

#### engProfile

One of the NVML\_COMPUTE\_INSTANCE\_ENGINE\_PROFILE\_\*

#### info

Returns detailed profile information

### Returns

- ▶ NVML\_SUCCESS Upon success
- ▶ NVML\_ERROR\_UNINITIALIZED If library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT If gpuInstance, profile, engProfile, info, or info->version are invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED If profile isn't supported
- ▶ NVML\_ERROR\_NO\_PERMISSION If user doesn't have permission to perform the operation

### Description

Versioned wrapper around `nvmlGpuInstanceGetComputeInstanceProfileInfo` that accepts a versioned `nvmlComputeInstanceProfileInfo_v2_t` or later output structure.



The caller must set the `nvmlComputeInstanceProfileInfo_v2_t::version` field to the appropriate version prior to calling this function. For example:

```
nvmlComputeInstanceProfileInfo_v2_t profileInfo =
    { .version = nvmlComputeInstanceProfileInfo_v2 };
nvmlReturn_t result
= nvmlGpuInstanceGetComputeInstanceProfileInfoV(gpuInstance,
    profile,
```

```
engProfile,
&profileInfo);
```

For Ampere or newer fully supported devices. Supported on Linux only.

## `nvmlReturn_t` `nvmlGpuInstanceGetComputeInstanceRemainingCapacity` (`nvmlGpuInstance_t` gpuInstance, unsigned int profileId, unsigned int \*count)

### Parameters

#### `gpuInstance`

The identifier of the target GPU instance

#### `profileId`

The compute instance profile ID. See

[nvmlGpuInstanceGetComputeInstanceProfileInfo](#)

#### `count`

Returns remaining instance count for the profile ID

### Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If `gpuInstance`, `profileId` or `availableCount` are invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` If `profileId` isn't supported
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation

### Description

Get compute instance profile capacity.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

## `nvmlReturn_t` `nvmlGpuInstanceGetComputeInstancePossiblePlacements` (`nvmlGpuInstance_t` gpuInstance, unsigned int profileId,

## `nvmlComputeInstancePlacement_t *placements,` `unsigned int *count)`

### Parameters

#### `gpuInstance`

The identifier of the target GPU instance

#### `profileId`

The compute instance profile ID. See

[nvmlGpuInstanceGetComputeInstanceProfileInfo](#)

#### `placements`

Returns placements allowed for the profile. Can be NULL to discover number of allowed placements for this profile. If non-NULL must be large enough to accommodate the placements supported by the profile.

#### `count`

Returns number of allowed placements for the profile.

### Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If `gpuInstance`, `profileId` or `count` are invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` If device doesn't have MIG mode enabled or `profileId` isn't supported
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation

### Description

Get compute instance placements.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

A placement represents the location of a compute instance within a GPU instance. This API only returns all the possible placements for the given profile. A created compute instance occupies compute slices described by its placement. Creation of new compute instance will fail if there is overlap with the already occupied compute slices.

**nvmlReturn\_t nvmlGpuInstanceCreateComputeInstance**  
 (nvmlGpuInstance\_t gpuInstance, unsigned int profileId,  
 nvmlComputeInstance\_t \*computeInstance)

#### Parameters

##### gpuInstance

The identifier of the target GPU instance

##### profileId

The compute instance profile ID. See

[nvmlGpuInstanceGetComputeInstanceProfileInfo](#)

##### computeInstance

Returns the compute instance handle

#### Returns

- ▶ NVML\_SUCCESS Upon success
- ▶ NVML\_ERROR\_UNINITIALIZED If library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT If gpuInstance, profile, profileId or computeInstance are invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED If profileId isn't supported
- ▶ NVML\_ERROR\_NO\_PERMISSION If user doesn't have permission to perform the operation
- ▶ NVML\_ERROR\_INSUFFICIENT\_RESOURCES If the requested compute instance could not be created

#### Description

Create compute instance.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

If the parent device is unbound, reset or the parent GPU instance is destroyed or the compute instance is destroyed explicitly, the compute instance handle would become invalid. The compute instance must be recreated to acquire a valid handle.

**nvmlReturn\_t**  
**nvmlGpuInstanceCreateComputeInstanceWithPlacement**  
 (nvmlGpuInstance\_t gpuInstance, unsigned int profileId,



```
const nvmlComputeInstancePlacement_t *placement,
nvmlComputeInstance_t *computeInstance)
```

### Parameters

#### **gpuInstance**

The identifier of the target GPU instance

#### **profileId**

The compute instance profile ID. See

[nvmlGpuInstanceGetComputeInstanceProfileInfo](#)

#### **placement**

The requested placement. See

[nvmlGpuInstanceGetComputeInstancePossiblePlacements](#)

#### **computeInstance**

Returns the compute instance handle

### Returns

- ▶ NVML\_SUCCESS Upon success
- ▶ NVML\_ERROR\_UNINITIALIZED If library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT If gpuInstance, profile, profileId or computeInstance are invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED If profileId isn't supported
- ▶ NVML\_ERROR\_NO\_PERMISSION If user doesn't have permission to perform the operation
- ▶ NVML\_ERROR\_INSUFFICIENT\_RESOURCES If the requested compute instance could not be created

### Description

Create compute instance with the specified placement.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

If the parent device is unbound, reset or the parent GPU instance is destroyed or the compute instance is destroyed explicitly, the compute instance handle would become invalid. The compute instance must be recreated to acquire a valid handle.

## `nvmlReturn_t nvmlComputeInstanceDestroy` (`nvmlComputeInstance_t computeInstance`)

### Parameters

#### `computeInstance`

The compute instance handle

### Returns

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If `computeInstance` is invalid
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation
- ▶ `NVML_ERROR_IN_USE` If the compute instance is in use. This error would be returned if processes (e.g. CUDA application) are active on the compute instance.

### Description

Destroy compute instance.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

## `nvmlReturn_t nvmlGpuInstanceGetComputeInstances` (`nvmlGpuInstance_t gpuInstance`, `unsigned int profileId`, `nvmlComputeInstance_t *computeInstances`, `unsigned int *count`)

### Parameters

#### `gpuInstance`

The identifier of the target GPU instance

#### `profileId`

The compute instance profile ID. See [nvmlGpuInstanceGetComputeInstanceProfileInfo](#)

#### `computeInstances`

Returns pre-existing compute instances, the buffer must be large enough to accommodate the instances supported by the profile. See [nvmlGpuInstanceGetComputeInstanceProfileInfo](#)

**count**

The count of returned compute instances

**Returns**

- ▶ NVML\_SUCCESS Upon success
- ▶ NVML\_ERROR\_UNINITIALIZED If library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT If gpuInstance, profileId, computeInstances or count are invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED If profileId isn't supported
- ▶ NVML\_ERROR\_NO\_PERMISSION If user doesn't have permission to perform the operation

**Description**

Get compute instances for given profile ID.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

**nvmlReturn\_t nvmlGpuInstanceGetComputeInstanceById (nvmlGpuInstance\_t gpuInstance, unsigned int id, nvmlComputeInstance\_t \*computeInstance)**

**Parameters****gpuInstance**

The identifier of the target GPU instance

**id**

The compute instance ID

**computeInstance**

Returns compute instance

**Returns**

- ▶ NVML\_SUCCESS Upon success
- ▶ NVML\_ERROR\_UNINITIALIZED If library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT If device, ID or computeInstance are invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED If device doesn't have MIG mode enabled
- ▶ NVML\_ERROR\_NO\_PERMISSION If user doesn't have permission to perform the operation
- ▶ NVML\_ERROR\_NOT\_FOUND If the compute instance is not found.

**Description**

Get compute instance for given instance ID.

For Ampere or newer fully supported devices. Supported on Linux only. Requires privileged user.

**`nvmlReturn_t nvmlComputeInstanceGetInfo_v2`**  
**`(nvmlComputeInstance_t computeInstance,`**  
**`nvmlComputeInstanceInfo_t *info)`**

**Parameters****computeInstance**

The compute instance handle

**info**

Return compute instance information

**Returns**

- ▶ `NVML_SUCCESS` Upon success
- ▶ `NVML_ERROR_UNINITIALIZED` If library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` If `computeInstance` or `info` are invalid
- ▶ `NVML_ERROR_NO_PERMISSION` If user doesn't have permission to perform the operation

**Description**

Get compute instance information.

For Ampere or newer fully supported devices. Supported on Linux only.

**`nvmlReturn_t nvmlDeviceIsMigDeviceHandle`**  
**`(nvmlDevice_t device, unsigned int *isMigDevice)`**

**Parameters****device**

NVML handle to test

**isMigDevice**

True when handle refers to a MIG device

**Returns**

- ▶ NVML\_SUCCESS if device status was successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device handle or isMigDevice reference is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this check is not supported by the device
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Test if the given handle refers to a MIG device.

A MIG device handle is an NVML abstraction which maps to a MIG compute instance. These overloaded references can be used (with some restrictions) interchangeably with a GPU device handle to execute queries at a per-compute instance granularity.

For Ampere or newer fully supported devices. Supported on Linux only.

## nvmlReturn\_t nvmlDeviceGetGpuInstanceId (nvmlDevice\_t device, unsigned int \*id)

**Parameters****device**

Target MIG device handle

**id**

GPU instance ID

**Returns**

- ▶ NVML\_SUCCESS if instance ID was successfully retrieved
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device or id reference is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Get GPU instance ID for the given MIG device handle.

GPU instance IDs are unique per device and remain valid until the GPU instance is destroyed.

For Ampere or newer fully supported devices. Supported on Linux only.

## `nvmlReturn_t nvmlDeviceGetComputeInstanceId` (`nvmlDevice_t device`, `unsigned int *id`)

### Parameters

**device**

Target MIG device handle

**id**

Compute instance ID

### Returns

- ▶ `NVML_SUCCESS` if instance ID was successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device or id reference is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Get compute instance ID for the given MIG device handle.

Compute instance IDs are unique per GPU instance and remain valid until the compute instance is destroyed.

For Ampere or newer fully supported devices. Supported on Linux only.

## `nvmlReturn_t nvmlDeviceGetMaxMigDeviceCount` (`nvmlDevice_t device`, `unsigned int *count`)

### Parameters

**device**

Target device handle

**count**

Count of MIG devices

### Returns

- ▶ `NVML_SUCCESS` if count was successfully retrieved
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if device or count reference is invalid
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

**Description**

Get the maximum number of MIG devices that can exist under a given parent NVML device.

Returns zero if MIG is not supported or enabled.

For Ampere or newer fully supported devices. Supported on Linux only.

**`nvmlReturn_t nvmlDeviceGetMigDeviceHandleByIndex`**  
**`(nvmlDevice_t device, unsigned int index, nvmlDevice_t`**  
**`*migDevice)`**

**Parameters****device**

Reference to the parent GPU device handle

**index**

Index of the MIG device

**migDevice**

Reference to the MIG device handle

**Returns**

- ▶ NVML\_SUCCESS if migDevice handle was successfully created
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device, index or migDevice reference is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device
- ▶ NVML\_ERROR\_NOT\_FOUND if no valid MIG device was found at index
- ▶ NVML\_ERROR\_UNKNOWN on any unexpected error

**Description**

Get MIG device handle for the given index under its parent NVML device.

If the compute instance is destroyed either explicitly or by destroying, resetting or unbinding the parent GPU instance or the GPU device itself the MIG device handle would remain invalid and must be requested again using this API. Handles may be reused and their properties can change in the process.

For Ampere or newer fully supported devices. Supported on Linux only.

## `nvmlReturn_t` `nvmlDeviceGetDeviceHandleFromMigDeviceHandle` (`nvmlDevice_t migDevice`, `nvmlDevice_t *device`)

### Parameters

#### **migDevice**

MIG device handle

#### **device**

Device handle

### Returns

- ▶ `NVML_SUCCESS` if device handle was successfully created
- ▶ `NVML_ERROR_UNINITIALIZED` if the library has not been successfully initialized
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if `migDevice` or `device` is invalid
- ▶ `NVML_ERROR_NOT_SUPPORTED` if this query is not supported by the device
- ▶ `NVML_ERROR_UNKNOWN` on any unexpected error

### Description

Get parent device handle from a MIG device handle.

For Ampere or newer fully supported devices. Supported on Linux only.

## `#define NVML_DEVICE_MIG_DISABLE 0x0`

Disable Multi Instance GPU mode.

## `#define NVML_DEVICE_MIG_ENABLE 0x1`

Enable Multi Instance GPU mode.

## `#define NVML_GPU_INSTANCE_PROFILE_1_SLICE 0x0`

GPU instance profiles.

These macros should be passed to `nvmlDeviceGetGpuInstanceProfileInfo` to retrieve the detailed information about a GPU instance such as profile ID, engine counts.

## `#define NVML_GPU_INSTANCE_PROFILE_CAPS_P2P 0x1`

MIG GPU instance profile capability.



Bit field values representing MIG profile capabilities  
 nvmlGpuInstanceProfileInfo\_v3\_t::capabilities

```
#define nvmlGpuInstanceProfileInfo_v2
NVML_STRUCT_VERSION(GpuInstanceProfileInfo, 2)
```

Version identifier value for nvmlGpuInstanceProfileInfo\_v2\_t::version.

```
#define nvmlGpuInstanceProfileInfo_v3
NVML_STRUCT_VERSION(GpuInstanceProfileInfo, 3)
```

Version identifier value for nvmlGpuInstanceProfileInfo\_v3\_t::version.

```
#define NVML_COMPUTE_INSTANCE_PROFILE_1_SLICE
0x0
```

Compute instance profiles.

These macros should be passed to `nvmlGpuInstanceGetComputeInstanceProfileInfo` to retrieve the detailed information about a compute instance such as profile ID, engine counts

```
#define
NVML_COMPUTE_INSTANCE_ENGINE_PROFILE_SHARED
0x0
```

All the engines except multiprocessors would be shared.

```
#define nvmlComputeInstanceProfileInfo_v2
NVML_STRUCT_VERSION(ComputeInstanceProfileInfo, 2)
```

Version identifier value for nvmlComputeInstanceProfileInfo\_v2\_t::version.

```
#define nvmlComputeInstanceProfileInfo_v3
NVML_STRUCT_VERSION(ComputeInstanceProfileInfo, 3)
```

Version identifier value for nvmlComputeInstanceProfileInfo\_v3\_t::version.

## 4.29. NVML GPM

## GPM Enums

## GPM Structs

## GPM Functions

### 4.29.1. GPM Enums

NVML GPM

#### enum nvmlGpmMetricId\_t

GPM Metric Identifiers

##### Values

**NVML\_GPM\_METRIC\_GRAPHICS\_UTIL = 1**

Percentage of time any compute/graphics app was active on the GPU. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_SM\_UTIL = 2**

Percentage of SMs that were busy. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_SM\_OCCUPANCY = 3**

Percentage of warps that were active vs theoretical maximum. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_INTEGER\_UTIL = 4**

Percentage of time the GPU's SMs were doing integer operations. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_ANY\_TENSOR\_UTIL = 5**

Percentage of time the GPU's SMs were doing ANY tensor operations. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_DFMA\_TENSOR\_UTIL = 6**

Percentage of time the GPU's SMs were doing DFMA tensor operations. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_HMMA\_TENSOR\_UTIL = 7**

Percentage of time the GPU's SMs were doing HMMA tensor operations. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_IMMA\_TENSOR\_UTIL = 9**

Percentage of time the GPU's SMs were doing IMMA tensor operations. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_DRAM\_BW\_UTIL = 10**

Percentage of DRAM bw used vs theoretical maximum. 0.0 - 100.0 %/.

**NVML\_GPM\_METRIC\_FP64\_UTIL = 11**

Percentage of time the GPU's SMs were doing non-tensor FP64 math. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_FP32\_UTIL = 12**

Percentage of time the GPU's SMs were doing non-tensor FP32 math. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_FP16\_UTIL = 13**

Percentage of time the GPU's SMs were doing non-tensor FP16 math. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_PCIE\_TX\_PER\_SEC = 20**

PCIe traffic from this GPU in MiB/sec.

**NVML\_GPM\_METRIC\_PCIE\_RX\_PER\_SEC = 21**  
PCIe traffic to this GPU in MiB/sec.

**NVML\_GPM\_METRIC\_NVDEC\_0\_UTIL = 30**  
Percent utilization of NVDEC 0. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVDEC\_1\_UTIL = 31**  
Percent utilization of NVDEC 1. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVDEC\_2\_UTIL = 32**  
Percent utilization of NVDEC 2. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVDEC\_3\_UTIL = 33**  
Percent utilization of NVDEC 3. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVDEC\_4\_UTIL = 34**  
Percent utilization of NVDEC 4. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVDEC\_5\_UTIL = 35**  
Percent utilization of NVDEC 5. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVDEC\_6\_UTIL = 36**  
Percent utilization of NVDEC 6. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVDEC\_7\_UTIL = 37**  
Percent utilization of NVDEC 7. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVJPG\_0\_UTIL = 40**  
Percent utilization of NVJPG 0. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVJPG\_1\_UTIL = 41**  
Percent utilization of NVJPG 1. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVJPG\_2\_UTIL = 42**  
Percent utilization of NVJPG 2. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVJPG\_3\_UTIL = 43**  
Percent utilization of NVJPG 3. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVJPG\_4\_UTIL = 44**  
Percent utilization of NVJPG 4. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVJPG\_5\_UTIL = 45**  
Percent utilization of NVJPG 5. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVJPG\_6\_UTIL = 46**  
Percent utilization of NVJPG 6. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVJPG\_7\_UTIL = 47**  
Percent utilization of NVJPG 7. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVOFA\_0\_UTIL = 50**  
Percent utilization of NVOFA 0. 0.0 - 100.0.

**NVML\_GPM\_METRIC\_NVLINK\_TOTAL\_RX\_PER\_SEC = 60**  
NvLink read bandwidth for all links in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_TOTAL\_TX\_PER\_SEC = 61**  
NvLink write bandwidth for all links in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L0\_RX\_PER\_SEC = 62**  
NvLink read bandwidth for link 0 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L0\_TX\_PER\_SEC = 63**

NvLink write bandwidth for link 0 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L1\_RX\_PER\_SEC = 64**

NvLink read bandwidth for link 1 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L1\_TX\_PER\_SEC = 65**

NvLink write bandwidth for link 1 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L2\_RX\_PER\_SEC = 66**

NvLink read bandwidth for link 2 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L2\_TX\_PER\_SEC = 67**

NvLink write bandwidth for link 2 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L3\_RX\_PER\_SEC = 68**

NvLink read bandwidth for link 3 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L3\_TX\_PER\_SEC = 69**

NvLink write bandwidth for link 3 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L4\_RX\_PER\_SEC = 70**

NvLink read bandwidth for link 4 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L4\_TX\_PER\_SEC = 71**

NvLink write bandwidth for link 4 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L5\_RX\_PER\_SEC = 72**

NvLink read bandwidth for link 5 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L5\_TX\_PER\_SEC = 73**

NvLink write bandwidth for link 5 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L6\_RX\_PER\_SEC = 74**

NvLink read bandwidth for link 6 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L6\_TX\_PER\_SEC = 75**

NvLink write bandwidth for link 6 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L7\_RX\_PER\_SEC = 76**

NvLink read bandwidth for link 7 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L7\_TX\_PER\_SEC = 77**

NvLink write bandwidth for link 7 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L8\_RX\_PER\_SEC = 78**

NvLink read bandwidth for link 8 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L8\_TX\_PER\_SEC = 79**

NvLink write bandwidth for link 8 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L9\_RX\_PER\_SEC = 80**

NvLink read bandwidth for link 9 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L9\_TX\_PER\_SEC = 81**

NvLink write bandwidth for link 9 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L10\_RX\_PER\_SEC = 82**

NvLink read bandwidth for link 10 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L10\_TX\_PER\_SEC = 83**

NvLink write bandwidth for link 10 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L11\_RX\_PER\_SEC = 84**

NvLink read bandwidth for link 11 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L11\_TX\_PER\_SEC = 85**

NvLink write bandwidth for link 11 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L12\_RX\_PER\_SEC = 86**

NvLink read bandwidth for link 12 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L12\_TX\_PER\_SEC = 87**

NvLink write bandwidth for link 12 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L13\_RX\_PER\_SEC = 88**

NvLink read bandwidth for link 13 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L13\_TX\_PER\_SEC = 89**

NvLink write bandwidth for link 13 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L14\_RX\_PER\_SEC = 90**

NvLink read bandwidth for link 14 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L14\_TX\_PER\_SEC = 91**

NvLink write bandwidth for link 14 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L15\_RX\_PER\_SEC = 92**

NvLink read bandwidth for link 15 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L15\_TX\_PER\_SEC = 93**

NvLink write bandwidth for link 15 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L16\_RX\_PER\_SEC = 94**

NvLink read bandwidth for link 16 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L16\_TX\_PER\_SEC = 95**

NvLink write bandwidth for link 16 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L17\_RX\_PER\_SEC = 96**

NvLink read bandwidth for link 17 in MiB/sec.

**NVML\_GPM\_METRIC\_NVLINK\_L17\_TX\_PER\_SEC = 97**

NvLink write bandwidth for link 17 in MiB/sec.

**NVML\_GPM\_METRIC\_MAX = 98**

Maximum value above +1. Note that changing this should also change

NVML\_GPM\_METRICS\_GET\_VERSION due to struct size change.

## 4.29.2. GPM Structs

NVML GPM

```
struct nvmlGpmMetric_t
```

```
struct nvmlGpmMetricsGet_t
```

```
struct nvmlGpmSupport_t
```

```
typedef struct nvmlGpmSample_st *nvmlGpmSample_t
```

Handle to an allocated GPM sample allocated with `nvmlGpmSampleAlloc()`. Free this with `nvmlGpmSampleFree()`.

### 4.29.3. GPM Functions

NVML GPM

```
nvmlReturn_t nvmlGpmMetricsGet (nvmlGpmMetricsGet_t
*metricsGet)
```

#### Parameters

**metricsGet**

IN/OUT: populated `nvmlGpmMetricsGet_t` struct

#### Returns

- ▶ NVML\_SUCCESS on success
- ▶ Nonzero NVML\_ERROR\_? enum on error

#### Description

Calculate GPM metrics from two samples.

For Hopper or newer fully supported devices.

```
nvmlReturn_t nvmlGpmSampleFree (nvmlGpmSample_t gpmSample)
```

#### Parameters

**gpmSample**

Sample to free

#### Returns

- ▶ NVML\_SUCCESS on success
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if an invalid pointer is provided

**Description**

Free an allocated sample buffer that was allocated with `nvmlGpmSampleAlloc()`

For Hopper or newer fully supported devices.

`nvmlReturn_t nvmlGpmSampleAlloc (nvmlGpmSample_t *gpmSample)`

**Parameters****gpmSample**

Where the allocated sample will be stored

**Returns**

- ▶ `NVML_SUCCESS` on success
- ▶ `NVML_ERROR_INVALID_ARGUMENT` if an invalid pointer is provided
- ▶ `NVML_ERROR_MEMORY` if system memory is insufficient

**Description**

Allocate a sample buffer to be used with NVML GPM . You will need to allocate at least two of these buffers to use with the NVML GPM feature

For Hopper or newer fully supported devices.

`nvmlReturn_t nvmlGpmSampleGet (nvmlDevice_t device, nvmlGpmSample_t gpmSample)`

**Parameters****device**

Device to get samples for

**gpmSample**

Buffer to read samples into

**Returns**

- ▶ `NVML_SUCCESS` on success
- ▶ Nonzero `NVML_ERROR_?` enum on error

**Description**

Read a sample of GPM metrics into the provided `gpmSample` buffer. After two samples are gathered, you can call `nvmlGpmMetricGet` on those samples to retrieve metrics

For Hopper or newer fully supported devices.

`nvmlReturn_t nvmlGpmMigSampleGet (nvmlDevice_t device, unsigned int gpuInstanceId, nvmlGpmSample_t gpmSample)`

### Parameters

#### **device**

Device to get samples for

#### **gpuInstanceId**

MIG GPU Instance ID

#### **gpmSample**

Buffer to read samples into

### Returns

- ▶ NVML\_SUCCESS on success
- ▶ Nonzero NVML\_ERROR\_? enum on error

### Description

Read a sample of GPM metrics into the provided gpmSample buffer for a MIG GPU Instance.

After two samples are gathered, you can call nvmlGpmMetricGet on those samples to retrieve metrics

For Hopper or newer fully supported devices.

`nvmlReturn_t nvmlGpmQueryDeviceSupport (nvmlDevice_t device, nvmlGpmSupport_t *gpmSupport)`

### Parameters

#### **device**

NVML device to query for

#### **gpmSupport**

Structure to indicate GPM support `nvmlGpmSupport_t`. Indicates GPM support per system for the supplied device

### Returns

- ▶ NVML\_SUCCESS on success
- ▶ Nonzero NVML\_ERROR\_? enum if there is an error in processing the query



**Description**

Indicate whether the supplied device supports GPM

**nvmlReturn\_t nvmlGpmQueryIfStreamingEnabled (nvmlDevice\_t device, unsigned int \*state)**

**Parameters****device**

The identifier of the target device

**state**

Returns GPM stream state NVML\_FEATURE\_DISABLED or NVML\_FEATURE\_ENABLED

**Returns**

- ▶ NVML\_SUCCESS if current GPM stream state were successfully queried
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid or state is NULL
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device

**Description**

Get GPM stream state.

HOPPER\_OR\_NEWER% Supported on Linux, Windows TCC.

**nvmlReturn\_t nvmlGpmSetStreamingEnabled (nvmlDevice\_t device, unsigned int state)**

**Parameters****device**

The identifier of the target device

**state**

GPM stream state, NVML\_FEATURE\_DISABLED or NVML\_FEATURE\_ENABLED

**Returns**

- ▶ NVML\_SUCCESS if current GPM stream state is successfully set
- ▶ NVML\_ERROR\_UNINITIALIZED if the library has not been successfully initialized
- ▶ NVML\_ERROR\_INVALID\_ARGUMENT if device is invalid
- ▶ NVML\_ERROR\_NOT\_SUPPORTED if this query is not supported by the device

**Description**

Set GPM stream state.

HOPPER\_OR\_NEWER% Supported on Linux, Windows TCC.

## 4.30. VirtualGPU

### vGPU Enums

### vGPU Constants

### vGPU Structs

#### 4.30.1. vGPU Enums

VirtualGPU

#### enum nvmlGpuVirtualizationMode\_t

GPU virtualization mode types.

**Values**

**NVML\_GPU\_VIRTUALIZATION\_MODE\_NONE = 0**

Represents Bare Metal GPU.

**NVML\_GPU\_VIRTUALIZATION\_MODE\_PASSTHROUGH = 1**

Device is associated with GPU-Passthrough.

**NVML\_GPU\_VIRTUALIZATION\_MODE\_VGPU = 2**

Device is associated with vGPU inside virtual machine.

**NVML\_GPU\_VIRTUALIZATION\_MODE\_HOST\_VGPU = 3**

Device is associated with VGX hypervisor in vGPU mode.

**NVML\_GPU\_VIRTUALIZATION\_MODE\_HOST\_VSGA = 4**

Device is associated with VGX hypervisor in vSGA mode.

#### enum nvmlHostVgpuMode\_t

Host vGPU modes

**Values**

**NVML\_HOST\_VGPU\_MODE\_NON\_SRIOV = 0**

Non SR-IOV mode.

**NVML\_HOST\_VGPU\_MODE\_SRIOV = 1**

SR-IOV mode.

## enum nvmlVgpuVmIdType\_t

Types of VM identifiers

### Values

**NVML\_VGPU\_VM\_ID\_DOMAIN\_ID = 0**

VM ID represents DOMAIN ID.

**NVML\_VGPU\_VM\_ID\_UUID = 1**

VM ID represents UUID.

## enum nvmlVgpuGuestInfoState\_t

vGPU GUEST info state

### Values

**NVML\_VGPU\_INSTANCE\_GUEST\_INFO\_STATE\_UNINITIALIZED = 0**

Guest-dependent fields uninitialized.

**NVML\_VGPU\_INSTANCE\_GUEST\_INFO\_STATE\_INITIALIZED = 1**

Guest-dependent fields initialized.

## enum nvmlGridLicenseFeatureCode\_t

vGPU software licensable features

### Values

**NVML\_GRID\_LICENSE\_FEATURE\_CODE\_UNKNOWN = 0**

Unknown.

**NVML\_GRID\_LICENSE\_FEATURE\_CODE\_VGPU = 1**

Virtual GPU.

**NVML\_GRID\_LICENSE\_FEATURE\_CODE\_NVIDIA\_RTX = 2**

Nvidia RTX.

**NVML\_GRID\_LICENSE\_FEATURE\_CODE\_VWORKSTATION =**

**NVML\_GRID\_LICENSE\_FEATURE\_CODE\_NVIDIA\_RTX**

Deprecated, do not use.

**NVML\_GRID\_LICENSE\_FEATURE\_CODE\_GAMING = 3**

Gaming.

**NVML\_GRID\_LICENSE\_FEATURE\_CODE\_COMPUTE = 4**

Compute.

## enum nvmlVgpuCapability\_t

vGPU queryable capabilities

### Values

**NVML\_VGPU\_CAP\_NVLINK\_P2P = 0**

P2P over NVLink is supported.

**NVML\_VGPU\_CAP\_GPUDIRECT = 1**

GPUDirect capability is supported.

**NVML\_VGPU\_CAP\_MULTI\_VGPU\_EXCLUSIVE = 2**

vGPU profile cannot be mixed with other vGPU profiles in same VM

**NVML\_VGPU\_CAP\_EXCLUSIVE\_TYPE = 3**

vGPU profile cannot run on a GPU alongside other profiles of different type

**NVML\_VGPU\_CAP\_EXCLUSIVE\_SIZE = 4**

vGPU profile cannot run on a GPU alongside other profiles of different size

**NVML\_VGPU\_CAP\_COUNT**

## enum nvmlVgpuDriverCapability\_t

vGPU driver queryable capabilities

### Values

**NVML\_VGPU\_DRIVER\_CAP\_HETEROGENEOUS\_MULTI\_VGPU = 0**

Supports mixing of different vGPU profiles within one guest VM.

**NVML\_VGPU\_DRIVER\_CAP\_COUNT**

## enum nvmlDeviceVgpuCapability\_t

Device vGPU queryable capabilities

### Values

**NVML\_DEVICE\_VGPU\_CAP\_FRACTIONAL\_MULTI\_VGPU = 0**

Query if the fractional vGPU profiles on this GPU can be used in multi-vGPU configurations.

**NVML\_DEVICE\_VGPU\_CAP\_HETEROGENEOUS\_TIMESLICE\_PROFILES = 1**

Query if the GPU support concurrent execution of timesliced vGPU profiles of differing types.

**NVML\_DEVICE\_VGPU\_CAP\_HETEROGENEOUS\_TIMESLICE\_SIZES = 2**

Query if the GPU support concurrent execution of timesliced vGPU profiles of differing framebuffer sizes.

**NVML\_DEVICE\_VGPU\_CAP\_READ\_DEVICE\_BUFFER\_BW = 3**

Query the GPU's read\_device\_buffer expected bandwidth capacity in megabytes per second.

**NVML\_DEVICE\_VGPU\_CAP\_WRITE\_DEVICE\_BUFFER\_BW = 4**

Query the GPU's write\_device\_buffer expected bandwidth capacity in megabytes per second.

**NVML\_DEVICE\_VGPU\_CAP\_DEVICE\_STREAMING = 5**

Query if vGPU profiles on the GPU supports migration data streaming.

**NVML\_DEVICE\_VGPU\_CAP\_MINI\_QUARTER\_GPU = 6**

Set/Get support for mini-quarter vGPU profiles.

**NVML\_DEVICE\_VGPU\_CAP\_COMPUTE\_MEDIA\_ENGINE\_GPU = 7**

Set/Get support for compute media engine vGPU profiles.

**NVML\_DEVICE\_VGPU\_CAP\_COUNT**

**#define NVML\_GRID\_LICENSE\_EXPIRY\_NOT\_AVAILABLE 0**

Expiry information not available.

Status codes for license expiry

**#define NVML\_GRID\_LICENSE\_EXPIRY\_INVALID 1**

Invalid expiry or error fetching expiry.

**#define NVML\_GRID\_LICENSE\_EXPIRY\_VALID 2**

Valid expiry.

**#define NVML\_GRID\_LICENSE\_EXPIRY\_NOT\_APPLICABLE 3**

Expiry not applicable.

**#define NVML\_GRID\_LICENSE\_EXPIRY\_PERMANENT 4**

Permanent expiry.

## 4.30.2. vGPU Constants

VirtualGPU

**#define NVML\_GRID\_LICENSE\_BUFFER\_SIZE 128**

Buffer size guaranteed to be large enough for `nvmlVgpuTypeGetLicense`

**#define NVML\_VGPU\_VIRTUALIZATION\_CAP\_MIGRATION 0:0**

Macros for vGPU instance's virtualization capabilities bitfield.

**#define NVML\_VGPU\_PGPU\_VIRTUALIZATION\_CAP\_MIGRATION 0:0**

Macros for pGPU's virtualization capabilities bitfield.

## 4.30.3. vGPU Structs

VirtualGPU

```
struct nvmlVgpuHeterogeneousMode_v1_t
struct nvmlVgpuPlacementId_v1_t
struct nvmlVgpuPlacementList_v1_t
struct nvmlVgpuInstanceUtilizationSample_t
struct nvmlVgpuInstanceUtilizationInfo_v1_t
struct nvmlVgpuInstancesUtilizationInfo_v1_t
struct nvmlVgpuProcessUtilizationSample_t
struct nvmlVgpuProcessUtilizationInfo_v1_t
struct nvmlVgpuProcessesUtilizationInfo_v1_t
union nvmlVgpuSchedulerParams_t
struct nvmlVgpuSchedulerLogEntry_t
struct nvmlVgpuSchedulerLog_t
struct nvmlVgpuSchedulerGetState_t
union nvmlVgpuSchedulerSetParams_t
struct nvmlVgpuSchedulerSetState_t
struct nvmlVgpuSchedulerCapabilities_t
struct nvmlVgpuLicenseExpiry_t
struct nvmlProcessUtilizationSample_t
struct nvmlProcessUtilizationInfo_v1_t
struct nvmlProcessesUtilizationInfo_v1_t
struct nvmlGridLicenseExpiry_t
```

```
struct nvmlGridLicensableFeature_t
struct nvmlGridLicensableFeatures_t
struct nvmlEccSramErrorStatus_v1_t
#define NVML_VGPU_SCHEDULER_POLICY_UNKNOWN 0
vGPU scheduler policies
#define NVML_GRID_LICENSE_STATE_UNKNOWN 0
Unknown state.
vGPU license state
#define NVML_GRID_LICENSE_STATE_UNINITIALIZED 1
Uninitialized state.
#define NVML_GRID_LICENSE_STATE_UNLICENSED_UNRESTRICTED 2
Unlicensed unrestricted state.
#define NVML_GRID_LICENSE_STATE_UNLICENSED_RESTRICTED 3
Unlicensed restricted state.
#define NVML_GRID_LICENSE_STATE_UNLICENSED 4
Unlicensed state.
#define NVML_GRID_LICENSE_STATE_LICENSED 5
Licensed state.
#define NVML_GSP_FIRMWARE_VERSION_BUF_SIZE 0x40
GSP firmware
#define NVML_DEVICE_ARCH_KEPLER 2
Simplified chip architecture
#define NVML_BUS_TYPE_UNKNOWN 0
PCI bus types
```



```
#define NVML_FAN_POLICY_TEMPERATURE_CONTINUOUS_SW 0
```

Device Power Modes Device Fan control policy

```
#define NVML_POWER_SOURCE_AC 0x00000000
```

Device Power Source

## 4.31. NvmlClocksEventReasons

```
#define nvmlClocksEventReasonGpuIdle  
0x00000000000000001LL
```

Nothing is running on the GPU and the clocks are dropping to Idle state



This limiter may be removed in a later release

```
#define  
nvmlClocksEventReasonApplicationsClocksSetting  
0x00000000000000002LL
```

GPU clocks are limited by current setting of applications clocks

**See also:**

[nvmlDeviceSetApplicationsClocks](#)

[nvmlDeviceGetApplicationsClock](#)

```
#define nvmlClocksThrottleReasonUserDefinedClocks  
nvmlClocksEventReasonApplicationsClocksSetting
```

Deprecated Renamed to [nvmlClocksThrottleReasonApplicationsClocksSetting](#) as the name describes the situation more accurately.

```
#define nvmlClocksEventReasonSwPowerCap  
0x00000000000000004LL
```

The clocks have been optimized to ensure not to exceed currently set power limits

**See also:**`nvidia-smi nvmlDeviceGetPowerUsage``nvidia-smi nvmlDeviceSetPowerManagementLimit``nvidia-smi nvmlDeviceGetPowerManagementLimit`**#define nvmlClocksThrottleReasonHwSlowdown  
0x000000000000000008LL**

HW Slowdown (reducing the core clocks by a factor of 2 or more) is engaged

This is an indicator of:

- ▶ temperature being too high
- ▶ External Power Brake Assertion is triggered (e.g. by the system power supply)
- ▶ Power draw is too high and Fast Trigger protection is reducing the clocks
- ▶ May be also reported during PState or clock change
  - ▶ This behavior may be removed in a later release.

**See also:**`nvidia-smi nvmlDeviceGetTemperature``nvidia-smi nvmlDeviceGetTemperatureThreshold``nvidia-smi nvmlDeviceGetPowerUsage`**#define nvmlClocksEventReasonSyncBoost  
0x000000000000000010LL**

Sync Boost

This GPU has been added to a Sync boost group with `nvidia-smi` or DCGM in order to maximize performance per watt. All GPUs in the sync boost group will boost to the minimum possible clocks across the entire group. Look at the throttle reasons for other GPUs in the system to see why those GPUs are holding this one at lower clocks.

**#define nvmlClocksEventReasonSwThermalSlowdown  
0x000000000000000020LL**

SW Thermal Slowdown

The current clocks have been optimized to ensure the the following is true:

- ▶ Current GPU temperature does not exceed GPU Max Operating Temperature

- ▶ Current memory temperature does not exceed Memory Max Operating Temperature

```
#define nvmClocksThrottleReasonHwThermalSlowdown
0x0000000000000040LL
```

HW Thermal Slowdown (reducing the core clocks by a factor of 2 or more) is engaged

This is an indicator of:

- ▶ temperature being too high

**See also:**

`nvmDeviceGetTemperature`

`nvmDeviceGetTemperatureThreshold`

`nvmDeviceGetPowerUsage`

```
#define
nvmClocksThrottleReasonHwPowerBrakeSlowdown
0x0000000000000080LL
```

HW Power Brake Slowdown (reducing the core clocks by a factor of 2 or more) is engaged

This is an indicator of:

- ▶ External Power Brake Assertion being triggered (e.g. by the system power supply)

**See also:**

`nvmDeviceGetTemperature`

`nvmDeviceGetTemperatureThreshold`

`nvmDeviceGetPowerUsage`

```
#define nvmClocksEventReasonDisplayClockSetting
0x0000000000000100LL
```

GPU clocks are limited by current setting of Display clocks

**See also:**

bug 1997531

```
#define nvmlClocksEventReasonNone
0x0000000000000000LL
```

Bit mask representing no clocks throttling

Clocks are as high as possible.

```
#define nvmlClocksEventReasonAll
(nvmlClocksThrottleReasonNone \ |
nvmlClocksEventReasonGpuIdle \ |
nvmlClocksEventReasonApplicationsClocksSetting
\ | nvmlClocksEventReasonSwPowerCap \
| nvmlClocksThrottleReasonHwSlowdown
\ | nvmlClocksEventReasonSyncBoost \ |
nvmlClocksEventReasonSwThermalSlowdown \ |
nvmlClocksThrottleReasonHwThermalSlowdown \ |
nvmlClocksThrottleReasonHwPowerBrakeSlowdown \ |
nvmlClocksEventReasonDisplayClockSetting \ )
```

Bit mask representing all supported clocks throttling reasons New reasons might be added to this list in the future

```
#define nvmlClocksThrottleReasonGpuIdle
nvmlClocksEventReasonGpuIdle
```

Deprecated Use `nvmlClocksEventReasonGpuIdle` instead

```
#define
nvmlClocksThrottleReasonApplicationsClocksSetting
nvmlClocksEventReasonApplicationsClocksSetting
```

Deprecated Use `nvmlClocksEventReasonApplicationsClocksSetting` instead

```
#define nvmlClocksThrottleReasonSyncBoost
nvmlClocksEventReasonSyncBoost
```

Deprecated Use `nvmlClocksEventReasonSyncBoost` instead

```
#define nvmlClocksThrottleReasonSwPowerCap  
nvmlClocksEventReasonSwPowerCap
```

Deprecated Use nvmlClocksEventReasonSwPowerCap instead

```
#define nvmlClocksThrottleReasonSwThermalSlowdown  
nvmlClocksEventReasonSwThermalSlowdown
```

Deprecated Use nvmlClocksEventReasonSwThermalSlowdown instead

```
#define nvmlClocksThrottleReasonDisplayClockSetting  
nvmlClocksEventReasonDisplayClockSetting
```

Deprecated Use nvmlClocksEventReasonDisplayClockSetting instead

```
#define nvmlClocksThrottleReasonNone  
nvmlClocksEventReasonNone
```

Deprecated Use nvmlClocksEventReasonNone instead

```
#define nvmlClocksThrottleReasonAll  
nvmlClocksEventReasonAll
```

Deprecated Use nvmlClocksEventReasonAll instead

# Chapter 5.

## DATA STRUCTURES

Here are the data structures with brief descriptions:

`nvmlAccountingStats_t`  
`nvmlBAR1Memory_t`  
`nvmlBridgeChipHierarchy_t`  
`nvmlBridgeChipInfo_t`  
`nvmlC2cModeInfo_v1_t`  
`nvmlClkMonFaultInfo_t`  
`nvmlClkMonStatus_t`  
`nvmlComputeInstanceProfileInfo_t`  
`nvmlComputeInstanceProfileInfo_v2_t`  
`nvmlComputeInstanceProfileInfo_v3_t`  
`nvmlConfComputeMemSizeInfo_t`  
`nvmlEccErrorCounts_t`  
`nvmlEccSramErrorStatus_v1_t`  
`nvmlEncoderSessionInfo_t`  
`nvmlEventData_t`  
`nvmlExcludedDeviceInfo_t`  
`nvmlFBCSessionInfo_t`  
`nvmlFBCStats_t`  
`nvmlFieldValue_t`  
`nvmlGpmMetric_t`  
`nvmlGpmMetricsGet_t`  
`nvmlGpmSupport_t`  
`nvmlGpuFabricInfo_v2_t`  
`nvmlGpuInstancePlacement_t`  
`nvmlGpuInstanceProfileInfo_t`  
`nvmlGpuInstanceProfileInfo_v2_t`  
`nvmlGpuInstanceProfileInfo_v3_t`  
`nvmlGridLicensableFeature_t`  
`nvmlGridLicensableFeatures_t`

`nvmlGridLicenseExpiry_t`  
`nvmlHwbcEntry_t`  
`nvmlLedState_t`  
`nvmlMemory_t`  
`nvmlMemory_v2_t`  
`nvmlNvLinkUtilizationControl_t`  
`nvmlPciInfo_t`  
`nvmlPciInfoExt_v1_t`  
`nvmlProcessDetail_v1_t`  
`nvmlProcessDetailList_v1_t`  
`nvmlProcessesUtilizationInfo_v1_t`  
`nvmlProcessInfo_t`  
`nvmlProcessInfo_v1_t`  
`nvmlProcessUtilizationInfo_v1_t`  
`nvmlProcessUtilizationSample_t`  
`nvmlPSUInfo_t`  
`nvmlRowRemapperHistogramValues_t`  
`nvmlSample_t`  
`nvmlSystemConfComputeSettings_v1_t`  
`nvmlUnitFanInfo_t`  
`nvmlUnitFanSpeeds_t`  
`nvmlUnitInfo_t`  
`nvmlUtilization_t`  
`nvmlValue_t`  
`nvmlVgpuHeterogeneousMode_v1_t`  
`nvmlVgpuInstancesUtilizationInfo_v1_t`  
`nvmlVgpuInstanceUtilizationInfo_v1_t`  
`nvmlVgpuInstanceUtilizationSample_t`  
`nvmlVgpuLicenseExpiry_t`  
`nvmlVgpuMetadata_t`  
`nvmlVgpuPgpuCompatibility_t`  
`nvmlVgpuPgpuMetadata_t`  
`nvmlVgpuPlacementId_v1_t`  
`nvmlVgpuPlacementList_v1_t`  
`nvmlVgpuProcessesUtilizationInfo_v1_t`  
`nvmlVgpuProcessUtilizationInfo_v1_t`  
`nvmlVgpuProcessUtilizationSample_t`  
`nvmlVgpuSchedulerCapabilities_t`  
`nvmlVgpuSchedulerGetState_t`  
`nvmlVgpuSchedulerLog_t`  
`nvmlVgpuSchedulerLogEntry_t`  
`nvmlVgpuSchedulerParams_t`  
`nvmlVgpuSchedulerSetParams_t`

`nvmlVgpuSchedulerSetState_t`  
`nvmlVgpuVersion_t`  
`nvmlViolationTime_t`

## 5.1. `nvmlAccountingStats_t` Struct Reference

Describes accounting statistics of a process.

### `unsigned int nvmlAccountingStats_t::gpuUtilization`

#### Description

Percent of time over the process's lifetime during which one or more kernels was executing on the GPU. Utilization stats just like returned by `nvmlDeviceGetUtilizationRates` but for the life time of a process (not just the last sample period). Set to `NVML_VALUE_NOT_AVAILABLE` if `nvmlDeviceGetUtilizationRates` is not supported

### `unsigned int nvmlAccountingStats_t::memoryUtilization`

#### Description

Percent of time over the process's lifetime during which global (device) memory was being read or written. Set to `NVML_VALUE_NOT_AVAILABLE` if `nvmlDeviceGetUtilizationRates` is not supported

### `unsigned long long nvmlAccountingStats_t::maxMemoryUsage`

#### Description

Maximum total memory in bytes that was ever allocated by the process. Set to `NVML_VALUE_NOT_AVAILABLE` if `nvmlProcessInfo_t->usedGpuMemory` is not supported

### `unsigned long long nvmlAccountingStats_t::time`

#### Description

Amount of time in ms during which the compute context was active. The time is reported as 0 if the process is not terminated



## **unsigned long long nvmlAccountingStats\_t::startTime**

CPU Timestamp in usec representing start time for the process.

## **unsigned int nvmlAccountingStats\_t::isRunning**

Flag to represent if the process is running (1 for running, 0 for terminated).

## **unsigned int nvmlAccountingStats\_t::reserved**

Reserved for future use.

## 5.2. nvmlBAR1Memory\_t Struct Reference

BAR1 Memory allocation Information for a device

### **unsigned long long nvmlBAR1Memory\_t::bar1Total**

Total BAR1 Memory (in bytes).

### **unsigned long long nvmlBAR1Memory\_t::bar1Free**

Unallocated BAR1 Memory (in bytes).

### **unsigned long long nvmlBAR1Memory\_t::bar1Used**

Allocated Used Memory (in bytes).

## 5.3. nvmlBridgeChipHierarchy\_t Struct Reference

This structure stores the complete Hierarchy of the Bridge Chip within the board. The immediate bridge is stored at index 0 of bridgeInfoList, parent to immediate bridge is at index 1 and so forth.

**unsigned char nvmlBridgeChipHierarchy\_t::bridgeCount**

Number of Bridge Chips on the Board.

**struct nvmlBridgeChipInfo\_t  
nvmlBridgeChipHierarchy\_t::bridgeChipInfo**

Hierarchy of Bridge Chips on the board.

## 5.4. nvmlBridgeChipInfo\_t Struct Reference

Information about the Bridge Chip Firmware

**nvmlBridgeChipType\_t nvmlBridgeChipInfo\_t::type**

Type of Bridge Chip.

**unsigned int nvmlBridgeChipInfo\_t::fwVersion**

Firmware Version. 0=Version is unavailable.

## 5.5. nvmlC2cModelInfo\_v1\_t Struct Reference

C2C Mode information for a device

## 5.6. nvmlClkMonFaultInfo\_t Struct Reference

Clock Monitor error types

**unsigned int nvmlClkMonFaultInfo\_t::clkApiDomain**

### Description

The Domain which faulted

**unsigned int  
nvmlClkMonFaultInfo\_t::clkDomainFaultMask**

### Description

Faults Information

## 5.7. nvmlClkMonStatus\_t Struct Reference

Clock Monitor Status

`unsigned int nvmlClkMonStatus_t::bGlobalStatus`

### Description

Fault status Indicator

`unsigned int nvmlClkMonStatus_t::clkMonListSize`

### Description

Total faulted domain numbers

`struct nvmlClkMonFaultInfo_t`  
`nvmlClkMonStatus_t::clkMonList`

### Description

The fault Information structure

## 5.8. nvmlComputeInstanceProfileInfo\_t Struct Reference

Compute instance profile information.

**unsigned int nvmlComputeInstanceProfileInfo\_t::id**

Unique profile ID within the GPU instance.

**unsigned int**

**nvmlComputeInstanceProfileInfo\_t::sliceCount**

GPU Slice count.

**unsigned int**

**nvmlComputeInstanceProfileInfo\_t::instanceCount**

Compute instance count.

**unsigned int**

**nvmlComputeInstanceProfileInfo\_t::multiprocessorCount**

Streaming Multiprocessor count.

**unsigned int**

**nvmlComputeInstanceProfileInfo\_t::sharedCopyEngineCount**

Shared Copy Engine count.

**unsigned int**

**nvmlComputeInstanceProfileInfo\_t::sharedDecoderCount**

Shared Decoder Engine count.

**unsigned int**

**nvmlComputeInstanceProfileInfo\_t::sharedEncoderCount**

Shared Encoder Engine count.

**unsigned int**

**nvmlComputeInstanceProfileInfo\_t::sharedJpegCount**

Shared JPEG Engine count.

**unsigned int**

**nvmlComputeInstanceProfileInfo\_t::sharedOfaCount**

Shared OFA Engine count.

## 5.9. nvmlComputeInstanceProfileInfo\_v2\_t Struct Reference

Compute instance profile information (v2).

Version 2 adds the `nvmlComputeInstanceProfileInfo_v2_t::version` field to the start of the structure, and the `nvmlComputeInstanceProfileInfo_v2_t::name` field to the end. This structure is not backwards-compatible with `nvmlComputeInstanceProfileInfo_t`.

unsigned int

`nvmlComputeInstanceProfileInfo_v2_t::version`

Structure version identifier (set to `nvmlComputeInstanceProfileInfo_v2`).

unsigned int `nvmlComputeInstanceProfileInfo_v2_t::id`

Unique profile ID within the GPU instance.

unsigned int

`nvmlComputeInstanceProfileInfo_v2_t::sliceCount`

GPU Slice count.

unsigned int

`nvmlComputeInstanceProfileInfo_v2_t::instanceCount`

Compute instance count.

unsigned int

`nvmlComputeInstanceProfileInfo_v2_t::multiprocessorCount`

Streaming Multiprocessor count.

unsigned int

`nvmlComputeInstanceProfileInfo_v2_t::sharedCopyEngineCount`

Shared Copy Engine count.

unsigned int

`nvmlComputeInstanceProfileInfo_v2_t::sharedDecoderCount`

Shared Decoder Engine count.

unsigned int

`nvmlComputeInstanceProfileInfo_v2_t::sharedEncoderCount`

Shared Encoder Engine count.

unsigned int

`nvmlComputeInstanceProfileInfo_v2_t::sharedJpegCount`

Shared JPEG Engine count.

unsigned int

`nvmlComputeInstanceProfileInfo_v2_t::sharedOfaCount`

Shared OFA Engine count.

`char nvmlComputeInstanceProfileInfo_v2_t::name`

Profile name.

## 5.10. `nvmlComputeInstanceProfileInfo_v3_t` Struct Reference

Compute instance profile information (v3).

Version 3 adds the `nvmlComputeInstanceProfileInfo_v3_t::capabilities` field

`nvmlComputeInstanceProfileInfo_t`.

unsigned int

`nvmlComputeInstanceProfileInfo_v3_t::version`

Structure version identifier (set to `nvmlComputeInstanceProfileInfo_v3`).

unsigned int `nvmlComputeInstanceProfileInfo_v3_t::id`

Unique profile ID within the GPU instance.

unsigned int

`nvmlComputeInstanceProfileInfo_v3_t::sliceCount`

GPU Slice count.

unsigned int

`nvmlComputeInstanceProfileInfo_v3_t::instanceCount`

Compute instance count.

unsigned int

`nvmlComputeInstanceProfileInfo_v3_t::multiprocessorCount`

Streaming Multiprocessor count.

unsigned int

`nvmlComputeInstanceProfileInfo_v3_t::sharedCopyEngineCount`

Shared Copy Engine count.

unsigned int

`nvmlComputeInstanceProfileInfo_v3_t::sharedDecoderCount`

Shared Decoder Engine count.

unsigned int

`nvmlComputeInstanceProfileInfo_v3_t::sharedEncoderCount`

Shared Encoder Engine count.

unsigned int

`nvmlComputeInstanceProfileInfo_v3_t::sharedJpegCount`

Shared JPEG Engine count.



**unsigned int**  
**nvmlComputeInstanceProfileInfo\_v3\_t::sharedOfaCount**  
Shared OFA Engine count.

**char nvmlComputeInstanceProfileInfo\_v3\_t::name**  
Profile name.

**unsigned int**  
**nvmlComputeInstanceProfileInfo\_v3\_t::capabilities**  
Additional capabilities.

## 5.11. nvmlConfComputeMemSizeInfo\_t Struct Reference

Protected memory size

## 5.12. nvmlEccErrorCounts\_t Struct Reference

Detailed ECC error counts for a device.

**Deprecated** Different GPU families can have different memory error counters See [nvmlDeviceGetMemoryErrorCounter](#)

`unsigned long long nvmlEccErrorCounts_t::l1Cache`

L1 cache errors.

`unsigned long long nvmlEccErrorCounts_t::l2Cache`

L2 cache errors.

`unsigned long long  
nvmlEccErrorCounts_t::deviceMemory`

Device memory errors.

`unsigned long long nvmlEccErrorCounts_t::registerFile`

Register file errors.

## 5.13. `nvmlEccSramErrorStatus_v1_t` Struct Reference

Structure to store SRAM uncorrectable error counters

**unsigned int nvmlEccSramErrorStatus\_v1\_t::version**

the API version number

**unsigned long long**

**nvmlEccSramErrorStatus\_v1\_t::aggregateUncParity**

aggregate uncorrectable parity error count

**unsigned long long**

**nvmlEccSramErrorStatus\_v1\_t::aggregateUncSecDed**

aggregate uncorrectable SEC-DED error count

**unsigned long long**

**nvmlEccSramErrorStatus\_v1\_t::aggregateCor**

aggregate correctable error count

**unsigned long long**

**nvmlEccSramErrorStatus\_v1\_t::volatileUncParity**

volatile uncorrectable parity error count

**unsigned long long**

**nvmlEccSramErrorStatus\_v1\_t::volatileUncSecDed**

volatile uncorrectable SEC-DED error count

**unsigned long long**

**nvmlEccSramErrorStatus\_v1\_t::volatileCor**

volatile correctable error count

**unsigned long long**

**nvmlEccSramErrorStatus\_v1\_t::aggregateUncBucketL2**

aggregate uncorrectable error count for L2 cache bucket

**unsigned long long**

**nvmlEccSramErrorStatus\_v1\_t::aggregateUncBucketSm**

aggregate uncorrectable error count for SM bucket

unsigned long long

`nvmlEccSramErrorStatus_v1_t::aggregateUncBucketPcie`

aggregate uncorrectable error count for PCIE bucket

unsigned long long

`nvmlEccSramErrorStatus_v1_t::aggregateUncBucketMcu`

aggregate uncorrectable error count for Microcontroller bucket

unsigned long long

`nvmlEccSramErrorStatus_v1_t::aggregateUncBucketOther`

aggregate uncorrectable error count for Other bucket

unsigned int

`nvmlEccSramErrorStatus_v1_t::bThresholdExceeded`

if the error threshold of field diag is exceeded

## 5.14. `nvmlEncoderSessionInfo_t` Struct Reference

Structure to hold encoder session data

**unsigned int nvmLEncoderSessionInfo\_t::sessionId**

Unique session ID.

**unsigned int nvmLEncoderSessionInfo\_t::pid**

Owning process ID.

**nvmLVgpuInstance\_t**

**nvmLEncoderSessionInfo\_t::vgpuInstance**

Owning vGPU instance ID (only valid on vGPU hosts, otherwise zero).

**nvmLEncoderType\_t**

**nvmLEncoderSessionInfo\_t::codecType**

Video encoder type.

**unsigned int nvmLEncoderSessionInfo\_t::hResolution**

Current encode horizontal resolution.

**unsigned int nvmLEncoderSessionInfo\_t::vResolution**

Current encode vertical resolution.

**unsigned int nvmLEncoderSessionInfo\_t::averageFps**

Moving average encode frames per second.

**unsigned int nvmLEncoderSessionInfo\_t::averageLatency**

Moving average encode latency in microseconds.

## 5.15. nvmLEventData\_t Struct Reference

Information about occurred event

**nvmlDevice\_t nvmlEventData\_t::device**

Specific device where the event occurred.

**unsigned long long nvmlEventData\_t::eventType**

Information about what specific event occurred.

**unsigned long long nvmlEventData\_t::eventId**

Stores XID error for the device in the event of nvmlEventTypeXidCriticalError.

**unsigned int nvmlEventData\_t::gpuInstanceId**

If MIG is enabled and nvmlEventTypeXidCriticalError event is attributable to a GPU.

**unsigned int nvmlEventData\_t::computeInstanceId**

If MIG is enabled and nvmlEventTypeXidCriticalError event is attributable to a.

## 5.16. nvmlExcludedDeviceInfo\_t Struct Reference

Excluded GPU device information

**struct nvmlPciInfo\_t nvmlExcludedDeviceInfo\_t::pciInfo**

The PCI information for the excluded GPU.

**char nvmlExcludedDeviceInfo\_t::uuid**

The ASCII string UUID for the excluded GPU.

## 5.17. nvmlFBCSessionInfo\_t Struct Reference

Structure to hold FBC session data

**unsigned int nvmlFBCSessionInfo\_t::sessionId**

Unique session ID.

**unsigned int nvmlFBCSessionInfo\_t::pid**

Owning process ID.

**nvmlVgpuInstance\_t**

**nvmlFBCSessionInfo\_t::vgpuInstance**

Owning vGPU instance ID (only valid on vGPU hosts, otherwise zero).

**unsigned int nvmlFBCSessionInfo\_t::displayOrdinal**

Display identifier.

**nvmlFBCSessionType\_t**

**nvmlFBCSessionInfo\_t::sessionType**

Type of frame buffer capture session.

**unsigned int nvmlFBCSessionInfo\_t::sessionFlags**

Session flags (one or more of NVML\_NVFBC\_SESSION\_FLAG\_XXX).

**unsigned int nvmlFBCSessionInfo\_t::hMaxResolution**

Max horizontal resolution supported by the capture session.

**unsigned int nvmlFBCSessionInfo\_t::vMaxResolution**

Max vertical resolution supported by the capture session.

**unsigned int nvmlFBCSessionInfo\_t::hResolution**

Horizontal resolution requested by caller in capture call.

**unsigned int nvmlFBCSessionInfo\_t::vResolution**

Vertical resolution requested by caller in capture call.

**unsigned int nvmlFBCSessionInfo\_t::averageFPS**

Moving average new frames captured per second.

**unsigned int nvmlFBCSessionInfo\_t::averageLatency**

Moving average new frame capture latency in microseconds.

## 5.18. nvmlFBCStats\_t Struct Reference

Structure to hold frame buffer capture sessions stats

**unsigned int nvmlFBCStats\_t::sessionsCount**

Total no of sessions.

**unsigned int nvmlFBCStats\_t::averageFPS**

Moving average new frames captured per second.

**unsigned int nvmlFBCStats\_t::averageLatency**

Moving average new frame capture latency in microseconds.

## 5.19. nvmlFieldValue\_t Struct Reference

Information for a Field Value Sample



## `unsigned int nvmlFieldValue_t::fieldId`

ID of the NVML field to retrieve. This must be set before any call that uses this struct. See the constants starting with `NVML_FI_` above.

## `unsigned int nvmlFieldValue_t::scopeId`

Scope ID can represent data used by NVML depending on `fieldId`'s context. For example, for NVLink throughput counter data, `scopeId` can represent `linkId`.

## `long long nvmlFieldValue_t::timestamp`

CPU Timestamp of this value in microseconds since 1970.

## `long long nvmlFieldValue_t::latencyUsec`

How long this field value took to update (in usec) within NVML. This may be averaged across several fields that are serviced by the same driver call.

## `nvmlValueType_t nvmlFieldValue_t::valueType`

Type of the value stored in `value`.

## `nvmlReturn_t nvmlFieldValue_t::nvmlReturn`

Return code for retrieving this value. This must be checked before looking at `value`, as `value` is undefined if `nvmlReturn != NVML_SUCCESS`.

## `nvmlFieldValue_t::value`

Value for this field. This is only valid if `nvmlReturn == NVML_SUCCESS`.

## 5.20. `nvmlGpmMetric_t` Struct Reference

GPM metric information.

**unsigned int nvmlGpmMetric\_t::metricId**

IN: NVML\_GPM\_METRIC\_? define of which metric to retrieve.

**nvmlReturn\_t nvmlGpmMetric\_t::nvmlReturn**

OUT: Status of this metric. If this is nonzero, then value is not valid.

**double nvmlGpmMetric\_t::value**

OUT: Value of this metric. Is only valid if nvmlReturn is 0 (NVML\_SUCCESS).

**nvmlGpmMetric\_t::@6 nvmlGpmMetric\_t::metricInfo**

OUT: Metric name and unit. Those can be NULL if not defined.

## 5.21. nvmlGpmMetricsGet\_t Struct Reference

GPM buffer information.

**unsigned int nvmlGpmMetricsGet\_t::version**

IN: Set to NVML\_GPM\_METRICS\_GET\_VERSION.

**unsigned int nvmlGpmMetricsGet\_t::numMetrics**

IN: How many metrics to retrieve in metrics[].

**nvmlGpmSample\_t nvmlGpmMetricsGet\_t::sample1**

IN: Sample buffer.

**nvmlGpmSample\_t nvmlGpmMetricsGet\_t::sample2**

IN: Sample buffer.

**struct nvmlGpmMetric\_t nvmlGpmMetricsGet\_t::metrics**

IN/OUT: Array of metrics. Set metricId on call. See nvmlReturn and value on return.

## 5.22. nvmlGpmSupport\_t Struct Reference

GPM device information.

## unsigned int nvmlGpmSupport\_t::version

IN: Set to NVML\_GPM\_SUPPORT\_VERSION.

## unsigned int nvmlGpmSupport\_t::isSupportedDevice

OUT: Indicates device support.

## 5.23. nvmlGpuFabricInfo\_v2\_t Struct Reference

GPU Fabric information (v2).

Version 2 adds the `nvmlGpuFabricInfo_v2_t::version` field to the start of the structure, and the `nvmlGpuFabricInfo_v2_t::healthMask` field to the end. This structure is not backwards-compatible with `nvmlGpuFabricInfo_t`.

### unsigned int nvmlGpuFabricInfo\_v2\_t::version

Structure version identifier (set to `nvmlGpuFabricInfo_v2`).

### unsigned char nvmlGpuFabricInfo\_v2\_t::clusterUuid

Uuid of the cluster to which this GPU belongs.

### nvmlReturn\_t nvmlGpuFabricInfo\_v2\_t::status

Error status, if any. Must be checked only if state returns "complete".

### unsigned int nvmlGpuFabricInfo\_v2\_t::cliQUEId

ID of the fabric clique to which this GPU belongs.

### nvmlGpuFabricState\_t nvmlGpuFabricInfo\_v2\_t::state

Current state of GPU registration process.

### unsigned int nvmlGpuFabricInfo\_v2\_t::healthMask

GPU Fabric health Status Mask.

## 5.24. nvmlGpuInstancePlacement\_t Struct Reference

MIG compute instance profile capability.

Bit field values representing MIG profile capabilities  
`nvmlComputeInstanceProfileInfo_v3_t::capabilities`

`unsigned int nvmlGpuInstancePlacement_t::start`

Index of first occupied memory slice.

`unsigned int nvmlGpuInstancePlacement_t::size`

Number of memory slices occupied.

## 5.25. `nvmlGpuInstanceProfileInfo_t` Struct Reference

GPU instance profile information.

**unsigned int nvmlGpuInstanceProfileInfo\_t::id**

Unique profile ID within the device.

**unsigned int  
nvmlGpuInstanceProfileInfo\_t::isP2pSupported**

Peer-to-Peer support.

**unsigned int nvmlGpuInstanceProfileInfo\_t::sliceCount**

GPU Slice count.

**unsigned int  
nvmlGpuInstanceProfileInfo\_t::instanceCount**

GPU instance count.

**unsigned int  
nvmlGpuInstanceProfileInfo\_t::multiprocessorCount**

Streaming Multiprocessor count.

**unsigned int  
nvmlGpuInstanceProfileInfo\_t::copyEngineCount**

Copy Engine count.

**unsigned int  
nvmlGpuInstanceProfileInfo\_t::decoderCount**

Decoder Engine count.

**unsigned int  
nvmlGpuInstanceProfileInfo\_t::encoderCount**

Encoder Engine count.

**unsigned int nvmlGpuInstanceProfileInfo\_t::jpegCount**

JPEG Engine count.

**unsigned int nvmlGpuInstanceProfileInfo\_t::ofaCount**

OFA Engine count.

unsigned long long  
`nvmlGpuInstanceProfileInfo_t::memorySizeMB`

Memory size in MBytes.

## 5.26. `nvmlGpuInstanceProfileInfo_v2_t` Struct Reference

GPU instance profile information (v2).

Version 2 adds the `nvmlGpuInstanceProfileInfo_v2_t::version` field to the start of the structure, and the `nvmlGpuInstanceProfileInfo_v2_t::name` field to the end. This structure is not backwards-compatible with `nvmlGpuInstanceProfileInfo_t`.

**unsigned int nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::version**

Structure version identifier (set to nvidia\_nvmlGpuInstanceProfileInfo\_v2).

**unsigned int nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::id**

Unique profile ID within the device.

**unsigned int**

**nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::isP2pSupported**

Peer-to-Peer support.

**unsigned int**

**nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::sliceCount**

GPU Slice count.

**unsigned int**

**nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::instanceCount**

GPU instance count.

**unsigned int**

**nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::multiprocessorCount**

Streaming Multiprocessor count.

**unsigned int**

**nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::copyEngineCount**

Copy Engine count.

**unsigned int**

**nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::decoderCount**

Decoder Engine count.

**unsigned int**

**nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::encoderCount**

Encoder Engine count.

**unsigned int**

**nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::jpegCount**

JPEG Engine count.

**unsigned int nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::ofaCount**

OFA Engine count.

**unsigned long long**

**nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::memorySizeMB**

Memory size in MBytes.

**char nvidia\_nvmlGpuInstanceProfileInfo\_v2\_t::name**

Profile name.

## 5.27. nvidia\_nvmlGpuInstanceProfileInfo\_v3\_t Struct Reference

GPU instance profile information (v3).

Version 3 removes `isP2pSupported` field and adds the `nvidia_nvmlGpuInstanceProfileInfo_v3_t::capabilities` field `nvidia_nvmlGpuInstanceProfileInfo_t`.



**unsigned int nvmGpuInstanceProfileInfo\_v3\_t::version**

Structure version identifier (set to nvmGpuInstanceProfileInfo\_v3).

**unsigned int nvmGpuInstanceProfileInfo\_v3\_t::id**

Unique profile ID within the device.

**unsigned int**

**nvmGpuInstanceProfileInfo\_v3\_t::sliceCount**

GPU Slice count.

**unsigned int**

**nvmGpuInstanceProfileInfo\_v3\_t::instanceCount**

GPU instance count.

**unsigned int**

**nvmGpuInstanceProfileInfo\_v3\_t::multiprocessorCount**

Streaming Multiprocessor count.

**unsigned int**

**nvmGpuInstanceProfileInfo\_v3\_t::copyEngineCount**

Copy Engine count.

**unsigned int**

**nvmGpuInstanceProfileInfo\_v3\_t::decoderCount**

Decoder Engine count.

**unsigned int**

**nvmGpuInstanceProfileInfo\_v3\_t::encoderCount**

Encoder Engine count.

**unsigned int**

**nvmGpuInstanceProfileInfo\_v3\_t::jpegCount**

JPEG Engine count.

**unsigned int nvmGpuInstanceProfileInfo\_v3\_t::ofaCount**

OFA Engine count.

**unsigned long long**  
**nvmlGpuInstanceProfileInfo\_v3\_t::memorySizeMB**

Memory size in MBytes.

**char nvmlGpuInstanceProfileInfo\_v3\_t::name**

Profile name.

**unsigned int**  
**nvmlGpuInstanceProfileInfo\_v3\_t::capabilities**

Additional capabilities.

## 5.28. nvmlGridLicensableFeature\_t Struct Reference

Structure containing vGPU software licensable feature information

`nvmlGridLicenseFeatureCode_t`  
`nvmlGridLicensableFeature_t::featureCode`

Licensed feature code.

`unsigned int nvmlGridLicensableFeature_t::featureState`

Non-zero if feature is currently licensed, otherwise zero.

`char nvmlGridLicensableFeature_t::licenseInfo`

Deprecated.

`char nvmlGridLicensableFeature_t::productName`

Product name of feature.

`unsigned int`

`nvmlGridLicensableFeature_t::featureEnabled`

Non-zero if feature is enabled, otherwise zero.

`struct nvmlGridLicenseExpiry_t`  
`nvmlGridLicensableFeature_t::licenseExpiry`

License expiry structure containing date and time.

## 5.29. `nvmlGridLicensableFeatures_t` Struct Reference

Structure to store vGPU software licensable features

**int**

**`nvmlGridLicensableFeatures_t::isGridLicenseSupported`**

Non-zero if vGPU Software Licensing is supported on the system, otherwise zero.

**unsigned int**

**`nvmlGridLicensableFeatures_t::licensableFeaturesCount`**

Entries returned in `gridLicensableFeatures` array.

**struct `nvmlGridLicensableFeature_t`**

**`nvmlGridLicensableFeatures_t::gridLicensableFeatures`**

Array of vGPU software licensable features.

## 5.30. `nvmlGridLicenseExpiry_t` Struct Reference

Structure to store license expiry date and time values

**unsigned int nvmlGridLicenseExpiry\_t::year**

Year value of license expiry.

**unsigned short nvmlGridLicenseExpiry\_t::month**

Month value of license expiry.

**unsigned short nvmlGridLicenseExpiry\_t::day**

Day value of license expiry.

**unsigned short nvmlGridLicenseExpiry\_t::hour**

Hour value of license expiry.

**unsigned short nvmlGridLicenseExpiry\_t::min**

Minutes value of license expiry.

**unsigned short nvmlGridLicenseExpiry\_t::sec**

Seconds value of license expiry.

**unsigned char nvmlGridLicenseExpiry\_t::status**

License expiry status.

## 5.31. nvmlHwbcEntry\_t Struct Reference

Description of HWBC entry

## 5.32. nvmlLedState\_t Struct Reference

LED states for an S-class unit.

## `char nvmlLedState_t::cause`

If amber, a text description of the cause.

## `nvmlLedColor_t nvmlLedState_t::color`

GREEN or AMBER.

## 5.33. `nvmlMemory_t` Struct Reference

Memory allocation information for a device (v1). The total amount is equal to the sum of the amounts of free and used memory.

### `unsigned long long nvmlMemory_t::total`

Total physical device memory (in bytes).

### `unsigned long long nvmlMemory_t::free`

Unallocated device memory (in bytes).

### `unsigned long long nvmlMemory_t::used`

#### **Description**

Sum of Reserved and Allocated device memory (in bytes). Note that the driver/GPU always sets aside a small amount of memory for bookkeeping

## 5.34. `nvmlMemory_v2_t` Struct Reference

Memory allocation information for a device (v2).

Version 2 adds versioning for the struct and the amount of system-reserved memory as an output.

### `unsigned int nvmlMemory_v2_t::version`

Structure format version (must be 2).

### `unsigned long long nvmlMemory_v2_t::total`

Total physical device memory (in bytes).

### `unsigned long long nvmlMemory_v2_t::reserved`

Device memory (in bytes) reserved for system use (driver or firmware).

### `unsigned long long nvmlMemory_v2_t::free`

Unallocated device memory (in bytes).

### `unsigned long long nvmlMemory_v2_t::used`

Allocated device memory (in bytes).

## 5.35. `nvmlNvLinkUtilizationControl_t` Struct Reference

Struct to define the NVLINK counter controls

## 5.36. `nvmlPciInfo_t` Struct Reference

PCI information about a GPU device.

### **char nvmnPciInfo\_t::busIdLegacy**

The legacy tuple domain:bus:device.function PCI identifier (& NULL terminator).

### **unsigned int nvmnPciInfo\_t::domain**

The PCI domain on which the device's bus resides, 0 to 0xffffffff.

### **unsigned int nvmnPciInfo\_t::bus**

The bus on which the device resides, 0 to 0xff.

### **unsigned int nvmnPciInfo\_t::device**

The device's id on the bus, 0 to 31.

### **unsigned int nvmnPciInfo\_t::pciDeviceId**

The combined 16-bit device id and 16-bit vendor id.

### **unsigned int nvmnPciInfo\_t::pciSubSystemId**

The 32-bit Sub System Device ID.

### **char nvmnPciInfo\_t::busId**

The tuple domain:bus:device.function PCI identifier (& NULL terminator).

## **5.37. nvmnPciInfoExt\_v1\_t Struct Reference**

PCI information about a GPU device.



**unsigned int nvmlPciInfoExt\_v1\_t::version**

The version number of this struct.

**unsigned int nvmlPciInfoExt\_v1\_t::domain**

The PCI domain on which the device's bus resides, 0 to 0xffffffff.

**unsigned int nvmlPciInfoExt\_v1\_t::bus**

The bus on which the device resides, 0 to 0xff.

**unsigned int nvmlPciInfoExt\_v1\_t::device**

The device's id on the bus, 0 to 31.

**unsigned int nvmlPciInfoExt\_v1\_t::pciDeviceId**

The combined 16-bit device id and 16-bit vendor id.

**unsigned int nvmlPciInfoExt\_v1\_t::pciSubSystemId**

The 32-bit Sub System Device ID.

**unsigned int nvmlPciInfoExt\_v1\_t::baseClass**

The 8-bit PCI base class code.

**unsigned int nvmlPciInfoExt\_v1\_t::subClass**

The 8-bit PCI sub class code.

**char nvmlPciInfoExt\_v1\_t::busId**

The tuple domain:bus:device.function PCI identifier (& NULL terminator).

## 5.38. nvmlProcessDetail\_v1\_t Struct Reference

Information about running process on the GPU with protected memory

**unsigned int nvidia\_nvmlProcessDetail\_v1\_t::pid**

Process ID.

**unsigned long long**

**nvidia\_nvmlProcessDetail\_v1\_t::usedGpuMemory**

#### Description

Amount of used GPU memory in bytes. Under WDDM, `NVML_VALUE_NOT_AVAILABLE` is always reported because Windows KMD manages all the memory and not the NVIDIA driver

**unsigned int nvidia\_nvmlProcessDetail\_v1\_t::gpuInstanceId**

If MIG is enabled, stores a valid GPU instance ID. `gpuInstanceId` is.

**unsigned int**

**nvidia\_nvmlProcessDetail\_v1\_t::computeInstanceId**

If MIG is enabled, stores a valid compute instance ID. `computeInstanceId`.

**unsigned long long**

**nvidia\_nvmlProcessDetail\_v1\_t::usedGpuCcProtectedMemory**

Amount of used GPU conf compute protected memory in bytes.

## 5.39. nvidia\_nvmlProcessDetailList\_v1\_t Struct Reference

Information about all running processes on the GPU for the given mode

**unsigned int nvmlProcessDetailList\_v1\_t::version**

Struct version, MUST be nvmlProcessDetailList\_v1.

**unsigned int nvmlProcessDetailList\_v1\_t::mode**

Process mode(Compute/Graphics/MPSCompute).

**unsigned int**

**nvmlProcessDetailList\_v1\_t::numProcArrayEntries**

Number of process entries in procArray.

**nvmlProcessDetail\_v1\_t**

**\*nvmlProcessDetailList\_v1\_t::procArray**

Process array.

## 5.40. nvmlProcessesUtilizationInfo\_v1\_t Struct Reference

Structure to store utilization and process ID for each running process -- version 1

**unsigned int nvmlProcessesUtilizationInfo\_v1\_t::version**

The version number of this struct.

**unsigned int**

**nvmlProcessesUtilizationInfo\_v1\_t::processSamplesCount**

Caller-supplied array size, and returns number of processes running.

**unsigned long long**

**nvmlProcessesUtilizationInfo\_v1\_t::lastSeenTimeStamp**

Return only samples with timestamp greater than lastSeenTimeStamp.

**nvmlProcessUtilizationInfo\_v1\_t**

**\*nvmlProcessesUtilizationInfo\_v1\_t::procUtilArray**

The array (allocated by caller) of the utilization of GPU SM, framebuffer, video encoder, video decoder, JPEG, and OFA.

## 5.41. nvmlProcessInfo\_t Struct Reference

Information about running compute processes on the GPU

**unsigned int nvmlProcessInfo\_t::pid**

Process ID.

**unsigned long long nvmlProcessInfo\_t::usedGpuMemory**

### Description

Amount of used GPU memory in bytes. Under WDDM, [NVML\\_VALUE\\_NOT\\_AVAILABLE](#) is always reported because Windows KMD manages all the memory and not the NVIDIA driver

## unsigned int nvmlProcessInfo\_t::gpuInstanceId

If MIG is enabled, stores a valid GPU instance ID. gpuInstanceId is set to.

## unsigned int nvmlProcessInfo\_t::computeInstanceId

If MIG is enabled, stores a valid compute instance ID. computeInstanceId is set to.

## 5.42. nvmlProcessInfo\_v1\_t Struct Reference

Information about running compute processes on the GPU, legacy version for older versions of the API.

### unsigned int nvmlProcessInfo\_v1\_t::pid

Process ID.

### unsigned long long nvmlProcessInfo\_v1\_t::usedGpuMemory

#### Description

Amount of used GPU memory in bytes. Under WDDM, `NVML_VALUE_NOT_AVAILABLE` is always reported because Windows KMD manages all the memory and not the NVIDIA driver

## 5.43. nvmlProcessUtilizationInfo\_v1\_t Struct Reference

Structure to store utilization value and process Id -- version 1

**unsigned long long**

**nvmlProcessUtilizationInfo\_v1\_t::timeStamp**

CPU Timestamp in microseconds.

**unsigned int nvmlProcessUtilizationInfo\_v1\_t::pid**

PID of process.

**unsigned int nvmlProcessUtilizationInfo\_v1\_t::smUtil**

SM (3D/Compute) Util Value.

**unsigned int nvmlProcessUtilizationInfo\_v1\_t::memUtil**

Frame Buffer Memory Util Value.

**unsigned int nvmlProcessUtilizationInfo\_v1\_t::encUtil**

Encoder Util Value.

**unsigned int nvmlProcessUtilizationInfo\_v1\_t::decUtil**

Decoder Util Value.

**unsigned int nvmlProcessUtilizationInfo\_v1\_t::jpgUtil**

Jpeg Util Value.

**unsigned int nvmlProcessUtilizationInfo\_v1\_t::ofaUtil**

Ofa Util Value.

## 5.44. nvmlProcessUtilizationSample\_t Struct Reference

Structure to store utilization value and process Id

**unsigned int nvmProcessUtilizationSample\_t::pid**

PID of process.

**unsigned long long**

**nvmProcessUtilizationSample\_t::timeStamp**

CPU Timestamp in microseconds.

**unsigned int nvmProcessUtilizationSample\_t::smUtil**

SM (3D/Compute) Util Value.

**unsigned int nvmProcessUtilizationSample\_t::memUtil**

Frame Buffer Memory Util Value.

**unsigned int nvmProcessUtilizationSample\_t::encUtil**

Encoder Util Value.

**unsigned int nvmProcessUtilizationSample\_t::decUtil**

Decoder Util Value.

## 5.45. nvmPSUInfo\_t Struct Reference

Power usage information for an S-class unit. The power supply state is a human readable string that equals "Normal" or contains a combination of "Abnormal" plus one or more of the following:

- ▶ High voltage
- ▶ Fan failure
- ▶ Heatsink temperature
- ▶ Current limit
- ▶ Voltage below UV alarm threshold
- ▶ Low-voltage
- ▶ SI2C remote off command
- ▶ MOD\_DISABLE input
- ▶ Short pin transition

### `char nvmI PSUInfo_t::state`

The power supply state.

### `unsigned int nvmI PSUInfo_t::current`

PSU current (A).

### `unsigned int nvmI PSUInfo_t::voltage`

PSU voltage (V).

### `unsigned int nvmI PSUInfo_t::power`

PSU power draw (W).

## 5.46. `nvmI RowRemapperHistogramValues_t` Struct Reference

Possible values that classify the remap availability for each bank. The max field will contain the number of banks that have maximum remap availability (all reserved rows are available). None means that there are no reserved rows available.

## 5.47. `nvmI Sample_t` Struct Reference

Information for Sample

### `unsigned long long nvmI Sample_t::timeStamp`

CPU Timestamp in microseconds.

### `nvmI Sample_t::sampleValue`

Sample Value.

## 5.48. `nvmI SystemConfComputeSettings_v1_t` Struct Reference

Confidential Compute System settings



## 5.49. nvmlUnitFanInfo\_t Struct Reference

Fan speed reading for a single fan in an S-class unit.

**unsigned int nvmlUnitFanInfo\_t::speed**

Fan speed (RPM).

**nvmlFanState\_t nvmlUnitFanInfo\_t::state**

Flag that indicates whether fan is working properly.

## 5.50. nvmlUnitFanSpeeds\_t Struct Reference

Fan speed readings for an entire S-class unit.

**struct nvmlUnitFanInfo\_t nvmlUnitFanSpeeds\_t::fans**

Fan speed data for each fan.

**unsigned int nvmlUnitFanSpeeds\_t::count**

Number of fans in unit.

## 5.51. nvmlUnitInfo\_t Struct Reference

Static S-class unit info.

## `char nvmlUnitInfo_t::name`

Product name.

## `char nvmlUnitInfo_t::id`

Product identifier.

## `char nvmlUnitInfo_t::serial`

Product serial number.

## `char nvmlUnitInfo_t::firmwareVersion`

Firmware version.

## 5.52. `nvmlUtilization_t` Struct Reference

Utilization information for a device. Each sample period may be between 1 second and 1/6 second, depending on the product being queried.

### `unsigned int nvmlUtilization_t::gpu`

Percent of time over the past sample period during which one or more kernels was executing on the GPU.

### `unsigned int nvmlUtilization_t::memory`

Percent of time over the past sample period during which global (device) memory was being read or written.

## 5.53. `nvmlValue_t` Union Reference

Union to represent different types of Value

**double nvmlValue\_t::dVal**

If the value is double.

**int nvmlValue\_t::siVal**

If the value is signed int.

**unsigned int nvmlValue\_t::uiVal**

If the value is unsigned int.

**unsigned long nvmlValue\_t::ulVal**

If the value is unsigned long.

**unsigned long long nvmlValue\_t::ullVal**

If the value is unsigned long long.

**signed long long nvmlValue\_t::sllVal**

If the value is signed long long.

## 5.54. nvmlVgpuHeterogeneousMode\_v1\_t Struct Reference

Structure to store the vGPU heterogeneous mode of device -- version 1

**unsigned int**

**nvmlVgpuHeterogeneousMode\_v1\_t::version**

The version number of this struct.

**unsigned int nvmlVgpuHeterogeneousMode\_v1\_t::mode**

The vGPU heterogeneous mode.

## 5.55. nvmlVgpuInstancesUtilizationInfo\_v1\_t Struct Reference

Structure to store recent utilization for vGPU instances running on a device -- version 1

**unsigned int**

**`nvmlVgpulInstancesUtilizationInfo_v1_t::version`**

The version number of this struct.

**`nvmlValueType_t`**

**`nvmlVgpulInstancesUtilizationInfo_v1_t::sampleValType`**

Hold the type of returned sample values.

**unsigned int**

**`nvmlVgpulInstancesUtilizationInfo_v1_t::vgpuInstanceCount`**

Hold the number of vGPU instances.

**unsigned long long**

**`nvmlVgpulInstancesUtilizationInfo_v1_t::lastSeenTimeStamp`**

Return only samples with timestamp greater than lastSeenTimeStamp.

**`nvmlVgpulInstanceUtilizationInfo_v1_t`**

**`*nvmlVgpulInstancesUtilizationInfo_v1_t::vgpuUtilArray`**

The array (allocated by caller) in which vGPU utilization are returned.

## 5.56. `nvmlVgpulInstanceUtilizationInfo_v1_t` Struct Reference

Structure to store Utilization Value and vgpuInstance Info -- Version 1

**unsigned long long**

**nvmlVgpulInstanceUtilizationInfo\_v1\_t::timeStamp**

CPU Timestamp in microseconds.

**nvmlVgpulInstance\_t**

**nvmlVgpulInstanceUtilizationInfo\_v1\_t::vgpulInstance**

vGPU Instance

**nvmlVgpulInstanceUtilizationInfo\_v1\_t::smUtil**

SM (3D/Compute) Util Value.

**nvmlVgpulInstanceUtilizationInfo\_v1\_t::memUtil**

Frame Buffer Memory Util Value.

**nvmlVgpulInstanceUtilizationInfo\_v1\_t::encUtil**

Encoder Util Value.

**nvmlVgpulInstanceUtilizationInfo\_v1\_t::decUtil**

Decoder Util Value.

**nvmlVgpulInstanceUtilizationInfo\_v1\_t::jpgUtil**

Jpeg Util Value.

**nvmlVgpulInstanceUtilizationInfo\_v1\_t::ofaUtil**

Ofa Util Value.

## 5.57. nvmlVgpulInstanceUtilizationSample\_t Struct Reference

Structure to store Utilization Value and vgpuInstance

**nvmlVgpulInstance\_t**  
**nvmlVgpulInstanceUtilizationSample\_t::vgpulInstance**  
vGPU Instance

**unsigned long long**  
**nvmlVgpulInstanceUtilizationSample\_t::timeStamp**  
CPU Timestamp in microseconds.

**nvmlVgpulInstanceUtilizationSample\_t::smUtil**  
SM (3D/Compute) Util Value.

**nvmlVgpulInstanceUtilizationSample\_t::memUtil**  
Frame Buffer Memory Util Value.

**nvmlVgpulInstanceUtilizationSample\_t::encUtil**  
Encoder Util Value.

**nvmlVgpulInstanceUtilizationSample\_t::decUtil**  
Decoder Util Value.

## 5.58. nvmlVgpuLicenseExpiry\_t Struct Reference

Structure to store the vGPU license expiry details

**unsigned int nvmIVgpuLicenseExpiry\_t::year**

Year of license expiry.

**unsigned short nvmIVgpuLicenseExpiry\_t::month**

Month of license expiry.

**unsigned short nvmIVgpuLicenseExpiry\_t::day**

Day of license expiry.

**unsigned short nvmIVgpuLicenseExpiry\_t::hour**

Hour of license expiry.

**unsigned short nvmIVgpuLicenseExpiry\_t::min**

Minutes of license expiry.

**unsigned short nvmIVgpuLicenseExpiry\_t::sec**

Seconds of license expiry.

**unsigned char nvmIVgpuLicenseExpiry\_t::status**

License expiry status.

## 5.59. nvmIVgpuMetadata\_t Struct Reference

vGPU metadata structure.

**unsigned int nvmiVgpuMetadata\_t::version**

Current version of the structure.

**unsigned int nvmiVgpuMetadata\_t::revision**

Current revision of the structure.

**nvmiVgpuGuestInfoState\_t**

**nvmiVgpuMetadata\_t::guestInfoState**

Current state of Guest-dependent fields.

**char nvmiVgpuMetadata\_t::guestDriverVersion**

Version of driver installed in guest.

**char nvmiVgpuMetadata\_t::hostDriverVersion**

Version of driver installed in host.

**unsigned int nvmiVgpuMetadata\_t::reserved**

Reserved for internal use.

**unsigned int**

**nvmiVgpuMetadata\_t::vgpuVirtualizationCaps**

vGPU virtualization capabilities bitfield

**unsigned int nvmiVgpuMetadata\_t::guestVgpuVersion**

vGPU version of guest driver

**unsigned int nvmiVgpuMetadata\_t::opaqueDataSize**

Size of opaque data field in bytes.

**char nvmiVgpuMetadata\_t::opaqueData**

Opaque data.

## 5.60. nvmiVgpuPgpuCompatibility\_t Struct Reference

vGPU-pGPU compatibility structure



**nvmlVgpuVmCompatibility\_t**

**nvmlVgpuPgpuCompatibility\_t::vgpuVmCompatibility**

Compatibility of vGPU VM. See `nvmlVgpuVmCompatibility_t`.

**nvmlVgpuPgpuCompatibilityLimitCode\_t**

**nvmlVgpuPgpuCompatibility\_t::compatibilityLimitCode**

Limiting factor for vGPU-pGPU compatibility. See `nvmlVgpuPgpuCompatibilityLimitCode_t`.

## 5.61. `nvmlVgpuPgpuMetadata_t` Struct Reference

Physical GPU metadata structure

**unsigned int nvmiVgpuPgpuMetadata\_t::version**

Current version of the structure.

**unsigned int nvmiVgpuPgpuMetadata\_t::revision**

Current revision of the structure.

**char nvmiVgpuPgpuMetadata\_t::hostDriverVersion**

Host driver version.

**unsigned int**

**nvmiVgpuPgpuMetadata\_t::pgpuVirtualizationCaps**

Pgpu virtualization capabilities bitfield.

**unsigned int nvmiVgpuPgpuMetadata\_t::reserved**

Reserved for internal use.

**struct nvmiVgpuVersion\_t**

**nvmiVgpuPgpuMetadata\_t::hostSupportedVgpuRange**

vGPU version range supported by host driver

**unsigned int nvmiVgpuPgpuMetadata\_t::opaqueDataSize**

Size of opaque data field in bytes.

**char nvmiVgpuPgpuMetadata\_t::opaqueData**

Opaque data.

## 5.62. nvmiVgpuPlacementId\_v1\_t Struct Reference

Structure to store the placement ID of vGPU instance -- version 1

**unsigned int nvmiVgpuPlacementId\_v1\_t::version**

The version number of this struct.

**unsigned int nvmiVgpuPlacementId\_v1\_t::placementId**

Placement ID of the active vGPU instance.

## 5.63. nvmiVgpuPlacementList\_v1\_t Struct Reference

Structure to store the list of vGPU placements -- version 1

**unsigned int nvmiVgpuPlacementList\_v1\_t::version**

The version number of this struct.

**unsigned int  
nvmiVgpuPlacementList\_v1\_t::placementSize**

The number of slots occupied by the vGPU type.

**unsigned int nvmiVgpuPlacementList\_v1\_t::count**

Count of placement IDs fetched.

**unsigned int  
\*nvmiVgpuPlacementList\_v1\_t::placementIds**

Placement IDs for the vGPU type.

## 5.64. nvmiVgpuProcessesUtilizationInfo\_v1\_t Struct Reference

Structure to store recent utilization, vgpuInstance and subprocess information for processes running on vGPU instances active on a device -- version 1

unsigned int

`nvmlVgpuProcessesUtilizationInfo_v1_t::version`

The version number of this struct.

unsigned int

`nvmlVgpuProcessesUtilizationInfo_v1_t::vgpuProcessCount`

Hold the number of processes running on vGPU instances.

unsigned long long

`nvmlVgpuProcessesUtilizationInfo_v1_t::lastSeenTimeStamp`

Return only samples with timestamp greater than lastSeenTimeStamp.

`nvmlVgpuProcessUtilizationInfo_v1_t`

`*nvmlVgpuProcessesUtilizationInfo_v1_t::vgpuProcUtilArray`

The array (allocated by caller) in which utilization of processes running on vGPU instances are returned.

## 5.65. `nvmlVgpuProcessUtilizationInfo_v1_t` Struct Reference

Structure to store Utilization Value, vgpuInstance and subprocess information for process running on vGPU instance -- version 1

**char nvmiVgpuProcessUtilizationInfo\_v1\_t::processName**

Name of process running within the vGPU VM.

**unsigned long long**

**nvmiVgpuProcessUtilizationInfo\_v1\_t::timeStamp**

CPU Timestamp in microseconds.

**nvmiVgpuInstance\_t**

**nvmiVgpuProcessUtilizationInfo\_v1\_t::vgpuInstance**

vGPU Instance

**unsigned int nvmiVgpuProcessUtilizationInfo\_v1\_t::pid**

PID of process running within the vGPU VM.

**unsigned int**

**nvmiVgpuProcessUtilizationInfo\_v1\_t::smUtil**

SM (3D/Compute) Util Value.

**unsigned int**

**nvmiVgpuProcessUtilizationInfo\_v1\_t::memUtil**

Frame Buffer Memory Util Value.

**unsigned int**

**nvmiVgpuProcessUtilizationInfo\_v1\_t::encUtil**

Encoder Util Value.

**unsigned int**

**nvmiVgpuProcessUtilizationInfo\_v1\_t::decUtil**

Decoder Util Value.

**unsigned int**

**nvmiVgpuProcessUtilizationInfo\_v1\_t::jpgUtil**

Jpeg Util Value.

**unsigned int**

**nvmiVgpuProcessUtilizationInfo\_v1\_t::ofaUtil**

Ofa Util Value.

## 5.66. nvmlVgpuProcessUtilizationSample\_t Struct Reference

Structure to store Utilization Value, vgpuInstance and subprocess information

`nvmlVgpuInstance_t`

`nvmlVgpuProcessUtilizationSample_t::vgpuInstance`

vGPU Instance

`unsigned int nvmlVgpuProcessUtilizationSample_t::pid`

PID of process running within the vGPU VM.

`char nvmlVgpuProcessUtilizationSample_t::processName`

Name of process running within the vGPU VM.

`unsigned long long`

`nvmlVgpuProcessUtilizationSample_t::timeStamp`

CPU Timestamp in microseconds.

`unsigned int`

`nvmlVgpuProcessUtilizationSample_t::smUtil`

SM (3D/Compute) Util Value.

`unsigned int`

`nvmlVgpuProcessUtilizationSample_t::memUtil`

Frame Buffer Memory Util Value.

`unsigned int`

`nvmlVgpuProcessUtilizationSample_t::encUtil`

Encoder Util Value.

`unsigned int`

`nvmlVgpuProcessUtilizationSample_t::decUtil`

Decoder Util Value.

## 5.67. `nvmlVgpuSchedulerCapabilities_t` Struct Reference

Structure to store the vGPU scheduler capabilities

unsigned int

`nvmlVgpuSchedulerCapabilities_t::supportedSchedulers`

List the supported vGPU schedulers on the device.

unsigned int

`nvmlVgpuSchedulerCapabilities_t::maxTimeslice`

Maximum timeslice value in ns.

unsigned int

`nvmlVgpuSchedulerCapabilities_t::minTimeslice`

Minimum timeslice value in ns.

unsigned int

`nvmlVgpuSchedulerCapabilities_t::isArrModeSupported`

Flag to check Adaptive Round Robin mode enabled/disabled.

unsigned int

`nvmlVgpuSchedulerCapabilities_t::maxFrequencyForARR`

Maximum frequency for Adaptive Round Robin mode.

unsigned int

`nvmlVgpuSchedulerCapabilities_t::minFrequencyForARR`

Minimum frequency for Adaptive Round Robin mode.

unsigned int

`nvmlVgpuSchedulerCapabilities_t::maxAvgFactorForARR`

Maximum averaging factor for Adaptive Round Robin mode.

unsigned int

`nvmlVgpuSchedulerCapabilities_t::minAvgFactorForARR`

Minimum averaging factor for Adaptive Round Robin mode.

## 5.68. `nvmlVgpuSchedulerGetState_t` Struct Reference

Structure to store the vGPU scheduler state



**unsigned int  
nvmiVgpuSchedulerGetState\_t::schedulerPolicy**

Scheduler policy.

**unsigned int nvmiVgpuSchedulerGetState\_t::arrMode**

Adaptive Round Robin scheduler mode. One of the NVML\_VGPU\_SCHEDULER\_ARR\_\*.

## 5.69. nvmiVgpuSchedulerLog\_t Struct Reference

Structure to store a vGPU software scheduler log

**unsigned int nvmiVgpuSchedulerLog\_t::engineId**

Engine whose software runlist log entries are fetched.

**unsigned int nvmiVgpuSchedulerLog\_t::schedulerPolicy**

Scheduler policy.

**unsigned int nvmiVgpuSchedulerLog\_t::arrMode**

Adaptive Round Robin scheduler mode. One of the NVML\_VGPU\_SCHEDULER\_ARR\_\*.

**unsigned int nvmiVgpuSchedulerLog\_t::entriesCount**

Count of log entries fetched.

## 5.70. nvmiVgpuSchedulerLogEntry\_t Struct Reference

Structure to store the state and logs of a software runlist

unsigned long long  
 nvmlVgpuSchedulerLogEntry\_t::timestamp

Timestamp in ns when this software runlist was preempted.

unsigned long long  
 nvmlVgpuSchedulerLogEntry\_t::timeRunTotal

Total time in ns this software runlist has run.

unsigned long long  
 nvmlVgpuSchedulerLogEntry\_t::timeRun

Time in ns this software runlist ran before preemption.

unsigned int nvmlVgpuSchedulerLogEntry\_t::swRunlistId

Software runlist Id.

unsigned long long  
 nvmlVgpuSchedulerLogEntry\_t::targetTimeSlice

The actual timeslice after deduction.

unsigned long long  
 nvmlVgpuSchedulerLogEntry\_t::cumulativePreemptionTime

Preemption time in ns for this SW runlist.

## 5.71. nvmlVgpuSchedulerParams\_t Union Reference

Union to represent the vGPU Scheduler Parameters

## `unsigned int nvmlVgpuSchedulerParams_t::avgFactor`

Average factor in compensating the timeslice for Adaptive Round Robin mode.

## `unsigned int nvmlVgpuSchedulerParams_t::timeslice`

The timeslice in ns for each software run list as configured, or the default value otherwise.

## 5.72. `nvmlVgpuSchedulerSetParams_t` Union Reference

Union to represent the vGPU Scheduler set Parameters

### `unsigned int nvmlVgpuSchedulerSetParams_t::avgFactor`

Average factor in compensating the timeslice for Adaptive Round Robin mode.

### `unsigned int nvmlVgpuSchedulerSetParams_t::frequency`

Frequency for Adaptive Round Robin mode.

### `unsigned int nvmlVgpuSchedulerSetParams_t::timeslice`

The timeslice in ns(Nanoseconds) for each software run list as configured, or the default value otherwise.

## 5.73. `nvmlVgpuSchedulerSetState_t` Struct Reference

Structure to set the vGPU scheduler state

unsigned int

`nvmlVgpuSchedulerSetState_t::schedulerPolicy`

Scheduler policy.

unsigned int

`nvmlVgpuSchedulerSetState_t::enableARRMode`

Adaptive Round Robin scheduler.

## 5.74. `nvmlVgpuVersion_t` Struct Reference

Structure representing range of vGPU versions.

unsigned int `nvmlVgpuVersion_t::minVersion`

Minimum vGPU version.

unsigned int `nvmlVgpuVersion_t::maxVersion`

Maximum vGPU version.

## 5.75. `nvmlViolationTime_t` Struct Reference

Struct to hold perf policy violation status data

unsigned long long `nvmlViolationTime_t::referenceTime`

`referenceTime` represents CPU timestamp in microseconds

unsigned long long `nvmlViolationTime_t::violationTime`

`violationTime` in Nanoseconds

# Chapter 6.

## DATA FIELDS

Here is a list of all documented struct and union fields with links to the struct/union documentation for each field:

### A

#### **aggregateCor**

[nvmlEccSramErrorStatus\\_v1\\_t](#)

#### **aggregateUncBucketL2**

[nvmlEccSramErrorStatus\\_v1\\_t](#)

#### **aggregateUncBucketMcu**

[nvmlEccSramErrorStatus\\_v1\\_t](#)

#### **aggregateUncBucketOther**

[nvmlEccSramErrorStatus\\_v1\\_t](#)

#### **aggregateUncBucketPcie**

[nvmlEccSramErrorStatus\\_v1\\_t](#)

#### **aggregateUncBucketSm**

[nvmlEccSramErrorStatus\\_v1\\_t](#)

#### **aggregateUncParity**

[nvmlEccSramErrorStatus\\_v1\\_t](#)

#### **aggregateUncSecDed**

[nvmlEccSramErrorStatus\\_v1\\_t](#)

#### **arrMode**

[nvmlVgpuSchedulerLog\\_t](#)

[nvmlVgpuSchedulerGetState\\_t](#)

#### **averageFps**

[nvmlEncoderSessionInfo\\_t](#)

#### **averageFPS**

[nvmlFBCSessionInfo\\_t](#)

[nvmlFBCStats\\_t](#)

#### **averageLatency**

[nvmlFBCStats\\_t](#)

nvmlFBCSessionInfo\_t  
 nvmlEncoderSessionInfo\_t

**avgFactor**

nvmlVgpuSchedulerParams\_t  
 nvmlVgpuSchedulerSetParams\_t

**B****bar1Free**

nvmlBAR1Memory\_t

**bar1Total**

nvmlBAR1Memory\_t

**bar1Used**

nvmlBAR1Memory\_t

**baseClass**

nvmlPciInfoExt\_v1\_t

**bGlobalStatus**

nvmlClkMonStatus\_t

**bridgeChipInfo**

nvmlBridgeChipHierarchy\_t

**bridgeCount**

nvmlBridgeChipHierarchy\_t

**bThresholdExceeded**

nvmlEccSramErrorStatus\_v1\_t

**bus**

nvmlPciInfoExt\_v1\_t  
 nvmlPciInfo\_t

**busId**

nvmlPciInfoExt\_v1\_t  
 nvmlPciInfo\_t

**busIdLegacy**

nvmlPciInfo\_t

**C****capabilities**

nvmlGpuInstanceProfileInfo\_v3\_t  
 nvmlComputeInstanceProfileInfo\_v3\_t

**cause**

nvmlLedState\_t

**cliqueId**

nvmlGpuFabricInfo\_v2\_t

**clkApiDomain**

nvmlClkMonFaultInfo\_t

**clkDomainFaultMask**

nvmClkMonFaultInfo\_t

**clkMonList**

nvmClkMonStatus\_t

**clkMonListSize**

nvmClkMonStatus\_t

**clusterUuid**

nvmGpuFabricInfo\_v2\_t

**codecType**

nvmEncoderSessionInfo\_t

**color**

nvmLedState\_t

**compatibilityLimitCode**

nvmVgpuPgpuCompatibility\_t

**computeInstanceId**

nvmEventData\_t

nvmProcessInfo\_t

nvmProcessDetail\_v1\_t

**copyEngineCount**

nvmGpuInstanceProfileInfo\_t

nvmGpuInstanceProfileInfo\_v3\_t

nvmGpuInstanceProfileInfo\_v2\_t

**count**

nvmUnitFanSpeeds\_t

nvmVgpuPlacementList\_v1\_t

**cumulativePreemptionTime**

nvmVgpuSchedulerLogEntry\_t

**current**

nvmPSUInfo\_t

**D****day**

nvmVgpuLicenseExpiry\_t

nvmGridLicenseExpiry\_t

**decoderCount**

nvmGpuInstanceProfileInfo\_v2\_t

nvmGpuInstanceProfileInfo\_v3\_t

nvmGpuInstanceProfileInfo\_t

**decUtil**

nvmVgpuProcessUtilizationSample\_t

nvmVgpuProcessUtilizationInfo\_v1\_t

nvmProcessUtilizationSample\_t

nvmProcessUtilizationInfo\_v1\_t

nvmlVgpuInstanceUtilizationSample\_t  
 nvmlVgpuInstanceUtilizationInfo\_v1\_t

**device**

nvmlPciInfo\_t  
 nvmlEventData\_t  
 nvmlPciInfoExt\_v1\_t

**deviceMemory**

nvmlEccErrorCounts\_t

**displayOrdinal**

nvmlFBCSessionInfo\_t

**domain**

nvmlPciInfoExt\_v1\_t  
 nvmlPciInfo\_t

**dVal**

nvmlValue\_t

**E****enableARRMode**

nvmlVgpuSchedulerSetState\_t

**encoderCount**

nvmlGpuInstanceProfileInfo\_t  
 nvmlGpuInstanceProfileInfo\_v3\_t  
 nvmlGpuInstanceProfileInfo\_v2\_t

**encUtil**

nvmlVgpuProcessUtilizationInfo\_v1\_t  
 nvmlProcessUtilizationInfo\_v1\_t  
 nvmlProcessUtilizationSample\_t  
 nvmlVgpuInstanceUtilizationSample\_t  
 nvmlVgpuProcessUtilizationSample\_t  
 nvmlVgpuInstanceUtilizationInfo\_v1\_t

**engineId**

nvmlVgpuSchedulerLog\_t

**entriesCount**

nvmlVgpuSchedulerLog\_t

**eventData**

nvmlEventData\_t

**eventType**

nvmlEventData\_t

**F****fans**

nvmlUnitFanSpeeds\_t



**featureCode**  
 nvmlGridLicensableFeature\_t  
**featureEnabled**  
 nvmlGridLicensableFeature\_t  
**featureState**  
 nvmlGridLicensableFeature\_t  
**fieldId**  
 nvmlFieldValue\_t  
**firmwareVersion**  
 nvmlUnitInfo\_t  
**free**  
 nvmlMemory\_t  
 nvmlMemory\_v2\_t  
**frequency**  
 nvmlVgpuSchedulerSetParams\_t  
**fwVersion**  
 nvmlBridgeChipInfo\_t

## G

**gpu**  
 nvmlUtilization\_t  
**gpuInstanceId**  
 nvmlProcessInfo\_t  
 nvmlEventData\_t  
 nvmlProcessDetail\_v1\_t  
**gpuUtilization**  
 nvmlAccountingStats\_t  
**gridLicensableFeatures**  
 nvmlGridLicensableFeatures\_t  
**guestDriverVersion**  
 nvmlVgpuMetadata\_t  
**guestInfoState**  
 nvmlVgpuMetadata\_t  
**guestVgpuVersion**  
 nvmlVgpuMetadata\_t

## H

**healthMask**  
 nvmlGpuFabricInfo\_v2\_t  
**hMaxResolution**  
 nvmlFBCSessionInfo\_t  
**hostDriverVersion**  
 nvmlVgpuPgpuMetadata\_t

`nvmlVgpuMetadata_t`

**hostSupportedVgpuRange**

`nvmlVgpuPgpuMetadata_t`

**hour**

`nvmlVgpuLicenseExpiry_t`

`nvmlGridLicenseExpiry_t`

**hResolution**

`nvmlFBCSessionInfo_t`

`nvmlEncoderSessionInfo_t`

**I**

**id**

`nvmlUnitInfo_t`

`nvmlGpuInstanceProfileInfo_t`

`nvmlGpuInstanceProfileInfo_v3_t`

`nvmlComputeInstanceProfileInfo_t`

`nvmlGpuInstanceProfileInfo_v2_t`

`nvmlComputeInstanceProfileInfo_v2_t`

`nvmlComputeInstanceProfileInfo_v3_t`

**instanceCount**

`nvmlGpuInstanceProfileInfo_t`

`nvmlGpuInstanceProfileInfo_v2_t`

`nvmlGpuInstanceProfileInfo_v3_t`

`nvmlComputeInstanceProfileInfo_t`

`nvmlComputeInstanceProfileInfo_v2_t`

`nvmlComputeInstanceProfileInfo_v3_t`

**isArrModeSupported**

`nvmlVgpuSchedulerCapabilities_t`

**isGridLicenseSupported**

`nvmlGridLicensableFeatures_t`

**isP2pSupported**

`nvmlGpuInstanceProfileInfo_t`

`nvmlGpuInstanceProfileInfo_v2_t`

**isRunning**

`nvmlAccountingStats_t`

**isSupportedDevice**

`nvmlGpmSupport_t`

**J**

**jpegCount**

`nvmlGpuInstanceProfileInfo_t`

`nvmlGpuInstanceProfileInfo_v2_t`

`nvmlGpuInstanceProfileInfo_v3_t`

**jpgUtil**

nvmlVgpuInstanceUtilizationInfo\_v1\_t  
 nvmlVgpuProcessUtilizationInfo\_v1\_t  
 nvmlProcessUtilizationInfo\_v1\_t

**L****l1Cache**

nvmlEccErrorCounts\_t

**l2Cache**

nvmlEccErrorCounts\_t

**lastSeenTimeStamp**

nvmlVgpuProcessesUtilizationInfo\_v1\_t  
 nvmlProcessesUtilizationInfo\_v1\_t  
 nvmlVgpuInstancesUtilizationInfo\_v1\_t

**latencyUsec**

nvmlFieldValue\_t

**licensableFeaturesCount**

nvmlGridLicensableFeatures\_t

**licenseExpiry**

nvmlGridLicensableFeature\_t

**licenseInfo**

nvmlGridLicensableFeature\_t

**M****maxAvgFactorForARR**

nvmlVgpuSchedulerCapabilities\_t

**maxFrequencyForARR**

nvmlVgpuSchedulerCapabilities\_t

**maxMemoryUsage**

nvmlAccountingStats\_t

**maxTimeslice**

nvmlVgpuSchedulerCapabilities\_t

**maxVersion**

nvmlVgpuVersion\_t

**memory**

nvmlUtilization\_t

**memorySizeMB**

nvmlGpuInstanceProfileInfo\_v2\_t  
 nvmlGpuInstanceProfileInfo\_v3\_t  
 nvmlGpuInstanceProfileInfo\_t

**memoryUtilization**

nvmlAccountingStats\_t

**memUtil**

nvmlProcessUtilizationInfo\_v1\_t  
 nvmlVgpuInstanceUtilizationSample\_t  
 nvmlVgpuInstanceUtilizationInfo\_v1\_t  
 nvmlVgpuProcessUtilizationSample\_t  
 nvmlVgpuProcessUtilizationInfo\_v1\_t  
 nvmlProcessUtilizationSample\_t

**metricId**

nvmlGpmMetric\_t

**metricInfo**

nvmlGpmMetric\_t

**metrics**

nvmlGpmMetricsGet\_t

**min**

nvmlVgpuLicenseExpiry\_t  
 nvmlGridLicenseExpiry\_t

**minAvgFactorForARR**

nvmlVgpuSchedulerCapabilities\_t

**minFrequencyForARR**

nvmlVgpuSchedulerCapabilities\_t

**minTimeslice**

nvmlVgpuSchedulerCapabilities\_t

**minVersion**

nvmlVgpuVersion\_t

**mode**

nvmlVgpuHeterogeneousMode\_v1\_t  
 nvmlProcessDetailList\_v1\_t

**month**

nvmlGridLicenseExpiry\_t  
 nvmlVgpuLicenseExpiry\_t

**multiprocessorCount**

nvmlGpuInstanceProfileInfo\_v3\_t  
 nvmlGpuInstanceProfileInfo\_v2\_t  
 nvmlComputeInstanceProfileInfo\_v3\_t  
 nvmlGpuInstanceProfileInfo\_t  
 nvmlComputeInstanceProfileInfo\_v2\_t  
 nvmlComputeInstanceProfileInfo\_t

**N****name**

nvmlUnitInfo\_t  
 nvmlGpuInstanceProfileInfo\_v2\_t  
 nvmlComputeInstanceProfileInfo\_v2\_t

`nvmlComputeInstanceProfileInfo_v3_t`

`nvmlGpuInstanceProfileInfo_v3_t`

### **numMetrics**

`nvmlGpmMetricsGet_t`

### **numProcArrayEntries**

`nvmlProcessDetailList_v1_t`

### **nvmlReturn**

`nvmlFieldValue_t`

`nvmlGpmMetric_t`

## **O**

### **ofaCount**

`nvmlGpuInstanceProfileInfo_t`

`nvmlGpuInstanceProfileInfo_v2_t`

`nvmlGpuInstanceProfileInfo_v3_t`

### **ofaUtil**

`nvmlVgpuInstanceUtilizationInfo_v1_t`

`nvmlVgpuProcessUtilizationInfo_v1_t`

`nvmlProcessUtilizationInfo_v1_t`

### **opaqueData**

`nvmlVgpuMetadata_t`

`nvmlVgpuPgpuMetadata_t`

### **opaqueDataSize**

`nvmlVgpuMetadata_t`

`nvmlVgpuPgpuMetadata_t`

## **P**

### **pciDeviceId**

`nvmlPciInfoExt_v1_t`

`nvmlPciInfo_t`

### **pciInfo**

`nvmlExcludedDeviceInfo_t`

### **pciSubSystemId**

`nvmlPciInfoExt_v1_t`

`nvmlPciInfo_t`

### **pgpuVirtualizationCaps**

`nvmlVgpuPgpuMetadata_t`

### **pid**

`nvmlProcessInfo_t`

`nvmlProcessDetail_v1_t`

`nvmlVgpuProcessUtilizationSample_t`

`nvmlVgpuProcessUtilizationInfo_v1_t`

`nvmlProcessUtilizationSample_t`

nvmlFBCSessionInfo\_t  
 nvmlEncoderSessionInfo\_t  
 nvmlProcessUtilizationInfo\_v1\_t  
 nvmlProcessInfo\_v1\_t

**placementId**

nvmlVgpuPlacementId\_v1\_t

**placementIds**

nvmlVgpuPlacementList\_v1\_t

**placementSize**

nvmlVgpuPlacementList\_v1\_t

**power**

nvmlPSUInfo\_t

**procArray**

nvmlProcessDetailList\_v1\_t

**processName**

nvmlVgpuProcessUtilizationInfo\_v1\_t

nvmlVgpuProcessUtilizationSample\_t

**processSamplesCount**

nvmlProcessesUtilizationInfo\_v1\_t

**procUtilArray**

nvmlProcessesUtilizationInfo\_v1\_t

**productName**

nvmlGridLicensableFeature\_t

**R****referenceTime**

nvmlViolationTime\_t

**registerFile**

nvmlEccErrorCounts\_t

**reserved**

nvmlAccountingStats\_t

nvmlVgpuPgpuMetadata\_t

nvmlVgpuMetadata\_t

nvmlMemory\_v2\_t

**revision**

nvmlVgpuPgpuMetadata\_t

nvmlVgpuMetadata\_t

**S****sample1**

nvmlGpmMetricsGet\_t

**sample2**

nvmlGpmMetricsGet\_t

**sampleValType**  
 nvmlVgpuInstancesUtilizationInfo\_v1\_t

**sampleValue**  
 nvmlSample\_t

**schedulerPolicy**  
 nvmlVgpuSchedulerLog\_t  
 nvmlVgpuSchedulerGetState\_t  
 nvmlVgpuSchedulerSetState\_t

**scopeId**  
 nvmlFieldValue\_t

**sec**  
 nvmlVgpuLicenseExpiry\_t  
 nvmlGridLicenseExpiry\_t

**serial**  
 nvmlUnitInfo\_t

**sessionFlags**  
 nvmlFBCSessionInfo\_t

**sessionId**  
 nvmlEncoderSessionInfo\_t  
 nvmlFBCSessionInfo\_t

**sessionsCount**  
 nvmlFBCStats\_t

**sessionType**  
 nvmlFBCSessionInfo\_t

**sharedCopyEngineCount**  
 nvmlComputeInstanceProfileInfo\_t  
 nvmlComputeInstanceProfileInfo\_v2\_t  
 nvmlComputeInstanceProfileInfo\_v3\_t

**sharedDecoderCount**  
 nvmlComputeInstanceProfileInfo\_t  
 nvmlComputeInstanceProfileInfo\_v2\_t  
 nvmlComputeInstanceProfileInfo\_v3\_t

**sharedEncoderCount**  
 nvmlComputeInstanceProfileInfo\_t  
 nvmlComputeInstanceProfileInfo\_v2\_t  
 nvmlComputeInstanceProfileInfo\_v3\_t

**sharedJpegCount**  
 nvmlComputeInstanceProfileInfo\_t  
 nvmlComputeInstanceProfileInfo\_v2\_t  
 nvmlComputeInstanceProfileInfo\_v3\_t

**sharedOfaCount**  
 nvmlComputeInstanceProfileInfo\_t  
 nvmlComputeInstanceProfileInfo\_v2\_t

nvmlComputeInstanceProfileInfo\_v3\_t  
**siVal**  
 nvmlValue\_t  
**size**  
 nvmlGpuInstancePlacement\_t  
**sliceCount**  
 nvmlGpuInstanceProfileInfo\_t  
 nvmlGpuInstanceProfileInfo\_v2\_t  
 nvmlGpuInstanceProfileInfo\_v3\_t  
 nvmlComputeInstanceProfileInfo\_t  
 nvmlComputeInstanceProfileInfo\_v2\_t  
 nvmlComputeInstanceProfileInfo\_v3\_t  
**sllVal**  
 nvmlValue\_t  
**smUtil**  
 nvmlVgpuInstanceUtilizationInfo\_v1\_t  
 nvmlVgpuInstanceUtilizationSample\_t  
 nvmlVgpuProcessUtilizationInfo\_v1\_t  
 nvmlProcessUtilizationSample\_t  
 nvmlVgpuProcessUtilizationSample\_t  
 nvmlProcessUtilizationInfo\_v1\_t  
**speed**  
 nvmlUnitFanInfo\_t  
**start**  
 nvmlGpuInstancePlacement\_t  
**startTime**  
 nvmlAccountingStats\_t  
**state**  
 nvmlUnitFanInfo\_t  
 nvmlGpuFabricInfo\_v2\_t  
 nvmlPSUInfo\_t  
**status**  
 nvmlGridLicenseExpiry\_t  
 nvmlGpuFabricInfo\_v2\_t  
 nvmlVgpuLicenseExpiry\_t  
**subClass**  
 nvmlPciInfoExt\_v1\_t  
**supportedSchedulers**  
 nvmlVgpuSchedulerCapabilities\_t  
**swRunlistId**  
 nvmlVgpuSchedulerLogEntry\_t



**T****targetTimeSlice**

`nvmlVgpuSchedulerLogEntry_t`

**time**

`nvmlAccountingStats_t`

**timeRun**

`nvmlVgpuSchedulerLogEntry_t`

**timeRunTotal**

`nvmlVgpuSchedulerLogEntry_t`

**timeslice**

`nvmlVgpuSchedulerParams_t`

`nvmlVgpuSchedulerSetParams_t`

**timeStamp**

`nvmlProcessUtilizationInfo_v1_t`

**timestamp**

`nvmlFieldValue_t`

**timeStamp**

`nvmlVgpuProcessUtilizationInfo_v1_t`

`nvmlVgpuInstanceUtilizationSample_t`

`nvmlVgpuInstanceUtilizationInfo_v1_t`

`nvmlVgpuProcessUtilizationSample_t`

`nvmlSample_t`

`nvmlProcessUtilizationSample_t`

**timestamp**

`nvmlVgpuSchedulerLogEntry_t`

**total**

`nvmlMemory_t`

`nvmlMemory_v2_t`

**type**

`nvmlBridgeChipInfo_t`

**U****uiVal**

`nvmlValue_t`

**ullVal**

`nvmlValue_t`

**ulVal**

`nvmlValue_t`

**used**

`nvmlMemory_t`

`nvmlMemory_v2_t`

**usedGpuCcProtectedMemory**

`nvmlProcessDetail_v1_t`

**usedGpuMemory**

nvmlProcessInfo\_v1\_t  
 nvmlProcessInfo\_t  
 nvmlProcessDetail\_v1\_t

**uuid**

nvmlExcludedDeviceInfo\_t

**V****value**

nvmlFieldValue\_t  
 nvmlGpmMetric\_t

**valueType**

nvmlFieldValue\_t

**version**

nvmlPciInfoExt\_v1\_t  
 nvmlVgpuPlacementId\_v1\_t  
 nvmlVgpuPgpuMetadata\_t  
 nvmlGpuInstanceProfileInfo\_v2\_t  
 nvmlVgpuPlacementList\_v1\_t  
 nvmlGpuInstanceProfileInfo\_v3\_t  
 nvmlComputeInstanceProfileInfo\_v2\_t  
 nvmlMemory\_v2\_t  
 nvmlVgpuInstancesUtilizationInfo\_v1\_t  
 nvmlComputeInstanceProfileInfo\_v3\_t  
 nvmlGpmMetricsGet\_t  
 nvmlVgpuProcessesUtilizationInfo\_v1\_t  
 nvmlGpmSupport\_t  
 nvmlProcessDetailList\_v1\_t  
 nvmlProcessesUtilizationInfo\_v1\_t  
 nvmlEccSramErrorStatus\_v1\_t  
 nvmlVgpuHeterogeneousMode\_v1\_t  
 nvmlGpuFabricInfo\_v2\_t  
 nvmlVgpuMetadata\_t

**vgpuInstance**

nvmlFBCSessionInfo\_t  
 nvmlVgpuInstanceUtilizationSample\_t  
 nvmlVgpuProcessUtilizationInfo\_v1\_t  
 nvmlVgpuInstanceUtilizationInfo\_v1\_t  
 nvmlEncoderSessionInfo\_t  
 nvmlVgpuProcessUtilizationSample\_t

**vgpuInstanceCount**

nvmlVgpuInstancesUtilizationInfo\_v1\_t

**vgpuProcessCount**  
    nvmlVgpuProcessesUtilizationInfo\_v1\_t

**vgpuProcUtilArray**  
    nvmlVgpuProcessesUtilizationInfo\_v1\_t

**vgpuUtilArray**  
    nvmlVgpuInstancesUtilizationInfo\_v1\_t

**vgpuVirtualizationCaps**  
    nvmlVgpuMetadata\_t

**vgpuVmCompatibility**  
    nvmlVgpuPgpuCompatibility\_t

**violationTime**  
    nvmlViolationTime\_t

**vMaxResolution**  
    nvmlFBCSessionInfo\_t

**volatileCor**  
    nvmlEccSramErrorStatus\_v1\_t

**volatileUncParity**  
    nvmlEccSramErrorStatus\_v1\_t

**volatileUncSecDed**  
    nvmlEccSramErrorStatus\_v1\_t

**voltage**  
    nvmlPSUInfo\_t

**vResolution**  
    nvmlEncoderSessionInfo\_t  
    nvmlFBCSessionInfo\_t

**Y**

**year**  
    nvmlVgpuLicenseExpiry\_t  
    nvmlGridLicenseExpiry\_t

# Chapter 7.

## DEPRECATED LIST

### **Class nvmlEccErrorCounts\_t**

Different GPU families can have different memory error counters See `nvmlDeviceGetMemoryErrorCounter`

### **Global nvmlEccBitType\_t**

See `nvmlMemoryErrorType_t` for a more flexible type

### **Global NVML\_SINGLE\_BIT\_ECC**

Mapped to `NVML_MEMORY_ERROR_TYPE_CORRECTED`

### **Global NVML\_DOUBLE\_BIT\_ECC**

Mapped to `NVML_MEMORY_ERROR_TYPE_UNCORRECTED`

### **Global nvmlDeviceGetHandleBySerial**

Since more than one GPU can exist on a single board this function is deprecated in favor of `nvmlDeviceGetHandleByUUID`. For dual GPU boards this function will return `NVML_ERROR_INVALID_ARGUMENT`.

### **Global nvmlDeviceGetCurrentClocksThrottleReasons**

Use `nvmlDeviceGetCurrentClocksEventReasons` instead

### **Global nvmlDeviceGetSupportedClocksThrottleReasons**

Use `nvmlDeviceGetSupportedClocksEventReasons` instead

**Global `nvmlDeviceGetDetailedEccErrors`**

This API supports only a fixed set of ECC error locations. On different GPU architectures different locations are supported. See `nvmlDeviceGetMemoryErrorCounter`.

**Global `nvmlVgpuInstanceGetLicenseStatus`**

Use `nvmlVgpuInstanceGetLicenseInfo_v2`.

**Global `nvmlClocksThrottleReasonUserDefinedClocks`**

Renamed to `nvmlClocksThrottleReasonApplicationsClocksSetting` as the name describes the situation more accurately.

**Global `nvmlClocksThrottleReasonGpuIdle`**

Use `nvmlClocksEventReasonGpuIdle` instead.

**Global `nvmlClocksThrottleReasonApplicationsClocksSetting`**

Use `nvmlClocksEventReasonApplicationsClocksSetting` instead.

**Global `nvmlClocksThrottleReasonSyncBoost`**

Use `nvmlClocksEventReasonSyncBoost` instead.

**Global `nvmlClocksThrottleReasonSwPowerCap`**

Use `nvmlClocksEventReasonSwPowerCap` instead.

**Global `nvmlClocksThrottleReasonSwThermalSlowdown`**

Use `nvmlClocksEventReasonSwThermalSlowdown` instead.

**Global `nvmlClocksThrottleReasonDisplayClockSetting`**

Use `nvmlClocksEventReasonDisplayClockSetting` instead.

**Global `nvmlClocksThrottleReasonNone`**

Use `nvmlClocksEventReasonNone` instead.

**Global nvmlClocksThrottleReasonAll**

Use nvmlClocksEventReasonAll instead

ALL NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, "MATERIALS") ARE BEING PROVIDED "AS IS." NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE.

Information furnished is believed to be accurate and reliable. However, NVIDIA Corporation assumes no responsibility for the consequences of use of such information or for any infringement of patents or other rights of third parties that may result from its use. No license is granted by implication of otherwise under any patent rights of NVIDIA Corporation. Specifications mentioned in this publication are subject to change without notice. This publication supersedes and replaces all other information previously supplied. NVIDIA Corporation products are not authorized as critical components in life support devices or systems without express written approval of NVIDIA Corporation.

NVIDIA and the NVIDIA logo are trademarks or registered trademarks of NVIDIA Corporation in the U.S. and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

© 2007-2024 NVIDIA Corporation & affiliates. All rights reserved.