



NVIDIA TensorRT

API Reference | NVIDIA Docs

Table of Contents

Chapter 1. Added, Deprecated, And Removed APIs.....	1
1.1. API Changes For TensorRT 8.4.0 EA.....	1
1.2. API Changes For TensorRT 8.2.0 EA.....	2
1.3. API Changes For TensorRT 8.0.1.....	3
Chapter 2. C++ API.....	11
Chapter 3. Python API.....	12

List of Tables

Table 1. New C++ APIs	1
Table 2. Deprecated C++ APIs	1
Table 3. New Python APIs	2
Table 4. Deprecated Python APIs	2
Table 5. New C++ APIs	2
Table 6. New Python APIs	3
Table 7. New C++ APIs	3
Table 8. Removed C++ APIs	4
Table 9. Removed Plugins	6
Table 10. Unsupported plugin methods removed in TensorRT 8.0	7
Table 11. Updated versions for supported plugin methods	7
Table 12. New Python APIs	7
Table 13. Removed Python APIs	8
Table 14. Deprecated APIs	10

Chapter 1. Added, Deprecated, And Removed APIs

1.1. API Changes For TensorRT 8.4.0 EA

The following tables show which APIs were added, deprecated, and removed for the NVIDIA® TensorRT™ 8.4.0 EA release.

C++ changes

Table 1. New C++ APIs

New C++ APIs
<u>setMemoryPoolLimit (IBuilderConfig::setMemoryPoolLimit)</u>
<u>getMemoryPoolLimit (IBuilderConfig::getMemoryPoolLimit)</u>
<u>MemoryPoolType</u>
<u>setMaxThreads (IBuilder::setMaxThreads, IRefitter::setMaxThreads, IRuntime::setMaxThreads)</u>
<u>getMaxThreads (IBuilder::getMaxThreads, IRefitter::getMaxThreads, IRuntime::getMaxThreads)</u>
<u>getBuilderPluginRegistry</u>

Table 2. Deprecated C++ APIs

Deprecated C++ APIs
<u>IFullyConnectedLayer</u>
<u>getMaxWorkspaceSize</u>
<u>setMaxWorkspaceSize</u>

Python changes

Table 3. New Python APIs

New Python APIs
<code>set_memory_pool_limit</code>
<code>get_memory_pool_limit</code>
<code>MemoryPoolType</code>
<code>max_threads</code> property (<code>Builder.max_threads</code> , <code>Refitter.max_threads</code> , <code>Runtime.max_threads</code>)
<code>get_builder_plugin_registry</code>

Table 4. Deprecated Python APIs

Deprecated Python APIs
<code>IFullyConnectedLayer</code>
<code>get_max_workspace_size</code>
<code>set_max_workspace_size</code>

1.2. API Changes For TensorRT 8.2.0 EA

The following tables show which APIs were added, deprecated, and removed for the NVIDIA[®] TensorRT™ 8.2.0 EA release.

C++ changes

Table 5. New C++ APIs

New C++ APIs
<code>IAssertionLayer</code>
<code>IConditionLayer</code>
<code>IEinsumLayer</code>
<code>IScatterLayer</code>

Python changes

Table 6. New Python APIs

New Python APIs
<u>IAssertionLayer</u>
<u>IConditionLayer</u>
<u>IEinsumLayer</u>
<u>IScatterLayer</u>

1.3. API Changes For TensorRT 8.0.1

The following tables show which APIs were added, deprecated, and removed for the TensorRT 8.0.1 release.

C++ changes

Table 7. New C++ APIs

New C++ APIs
<u>class IDequantizeLayer</u>
<u>class IQuantizeLayer</u>
<u>class ITimingCache</u>
<u>IBuilder::buildSerializedNetwork()</u>
<u>IBuilderConfig::getTimingCache()</u>
<u>IBuilderConfig::setTimingCache()</u>
<u>IGpuAllocator::reallocate()</u>
<u>INetworkDefinition::addDequantize()</u>
<u>INetworkDefinition::addQuantize()</u>
<u>INetworkDefinition::setWeightsName()</u>
<u>IPluginRegistry::deregisterCreator()</u>
<u>IRefitter::getMissingWeights()</u>
<u>IRefitter::getAllWeights()</u>
<u>IRefitter::setNamedWeights()</u>
<u>IResizeLayer::getCoordinateTransformation()</u>

New C++ APIs
<u>IResizeLayer::getNearestRounding()</u>
<u>IResizeLayer::getSelectorForSinglePixel()</u>
<u>IResizeLayer::setCoordinateTransformation()</u>
<u>IResizeLayer::setNearestRounding()</u>
<u>IResizeLayer::setSelectorForSinglePixel()</u>
<u>IScaleLayer::setChannelAxis()</u>
<u>enum ResizeCoordinateTransformation</u>
<u>enum ResizeMode</u>
<u>BuilderFlag::kSPARSE_WEIGHTS</u>
<u>TacticSource::kCUDNN</u>
<u>TensorFormat::kDLA_HWC4</u>
<u>TensorFormat::kDLA_LINEAR</u>
<u>TensorFormat::kHWC16</u>

Table 8. Removed C++ APIs

Removed C++ APIs
Core Library
DimensionType
Dims::Type
class DimsCHW
class DimsNCHW
class IOutputDimensionFormula
class IPlugin
class IPluginFactory
class IPluginLayer
class IRNNLayer
<u>IBuilder::getEngineCapability()</u>
<u>IBuilder::allowGPUFallback()</u>
<u>IBuilder::buildCudaEngine()</u>
<u>IBuilder::canRunOnDLA()</u>
<u>IBuilder::createNetwork()</u>
<u>IBuilder::getAverageFindIterations()</u>

Removed C++ APIs

<code>IBuilder::getDebugSync()</code>
<code>IBuilder::getDefaultDeviceType()</code>
<code>IBuilder::getDeviceType()</code>
<code>IBuilder::getDLACore()</code>
<code>IBuilder::getFp16Mode()</code>
<code>IBuilder::getHalf2Mode()</code>
<code>IBuilder::getInt8Mode()</code>
<code>IBuilder::getMaxWorkspaceSize()</code>
<code>IBuilder::getMinFindIterations()</code>
<code>IBuilder::getRefittable()</code>
<code>IBuilder::getStrictTypeConstraints()</code>
<code>IBuilder::isDeviceTypeSet()</code>
<code>IBuilder::reset()</code>
<code>IBuilder::resetDeviceType()</code>
<code>IBuilder::setAverageFindIterations()</code>
<code>IBuilder::setDebugSync()</code>
<code>IBuilder::setDefaultDeviceType()</code>
<code>IBuilder::setDeviceType()</code>
<code>IBuilder::setDLACore()</code>
<code>IBuilder::setEngineCapability()</code>
<code>IBuilder::setFp16Mode()</code>
<code>IBuilder::setHalf2Mode()</code>
<code>IBuilder::setInt8Calibrator()</code>
<code>IBuilder::setInt8Mode()</code>
<code>IBuilder::setMaxWorkspaceSize()</code>
<code>IBuilder::setMinFindIterations()</code>
<code>IBuilder::setRefittable()</code>
<code>IBuilder::setStrictTypeConstraints()</code>
<code>ICudaEngine::getWorkspaceSize()</code>
<code>IMatrixMultiplyLayer::getTranspose()</code>
<code>IMatrixMultiplyLayer::setTranspose()</code>
<code>INetworkDefinition::addMatrixMultiply()</code>

Removed C++ APIs
<code>INetworkDefinition::addPlugin()</code>
<code>INetworkDefinition::addPluginExt()</code>
<code>INetworkDefinition::addRNN()</code>
<code>INetworkDefinition::getConvolutionOutputDimensionsFormula()</code>
<code>INetworkDefinition::getDeconvolutionOutputDimensionsFormula()</code>
<code>INetworkDefinition::getPoolingOutputDimensionsFormula()</code>
<code>INetworkDefinition::setConvolutionOutputDimensionsFormula()</code>
<code>INetworkDefinition::setDeconvolutionOutputDimensionsFormula()</code>
<code>INetworkDefinition::setPoolingOutputDimensionsFormula()</code>
<code>ITensor::getDynamicRange()</code>
<code>TensorFormat::kNHWC8</code>
<code>TensorFormat::kNCHW</code>
<code>TensorFormat::kNC2HW2</code>
Caffe Parser
<code>class IPluginFactory</code>
<code>class IPluginFactoryExt</code>
<code>setPluginFactory()</code>
<code>setPluginFactoryExt()</code>
UFF Parser
<code>class IPluginFactory</code>
<code>class IPluginFactoryExt</code>
<code>setPluginFactory()</code>
<code>setPluginFactoryExt()</code>

Table 9. Removed Plugins

Removed Plugins
<code>class INvPlugin</code>
<code>createLReLUPlugin()</code>
<code>createClipPlugin()</code>
<code>PluginType</code>
<code>struct SoftmaxTree</code>

For plugins based on `IPluginV2DynamicExt` and `IPluginV2IOExt`, certain methods with legacy function signatures (derived from `IPluginV2` and `IPluginV2Ext` base classes) which were deprecated and marked for removal in TensorRT 8.0 will no longer be available. Plugins using these interface methods must stop using them or implement the versions with updated signatures, as applicable.

Table 10. Unsupported plugin methods removed in TensorRT 8.0

Removed Plugins
<code>IPluginV2DynamicExt::canBroadcastInputAcrossBatch()</code>
<code>IPluginV2DynamicExt::isOutputBroadcastAcrossBatch()</code>
<code>IPluginV2DynamicExt::getTensorRTVersion()</code>
<code>IPluginV2IOExt::configureWithFormat()</code>
<code>IPluginV2IOExt::getTensorRTVersion()</code>

Table 11. Updated versions for supported plugin methods

Removed Plugin	Replaced with
	<code>IPluginV2DynamicExt::configurePlugin()</code>
	<code>IPluginV2DynamicExt::enqueue()</code>
	<code>IPluginV2DynamicExt::getOutputDimensions()</code>
	<code>IPluginV2DynamicExt::getWorkspaceSize()</code>
	<code>IPluginV2IOExt::configurePlugin()</code>
<code>IPluginV2DynamicExt::supportsFormat()</code>	<code>IPluginV2DynamicExt::supportsFormatCombination()</code>
<code>IPluginV2IOExt::supportsFormat()</code>	<code>IPluginV2IOExt::supportsFormatCombination()</code>

Python changes

Table 12. New Python APIs

New Python APIs
<code>class IDequantizeLayer</code>
<code>class IQuantizeLayer</code>
<code>class ITimingCache</code>
<code>Builder.build_serialized_network()</code>
<code>IBuilderConfig.get_timing_cache()</code>

New Python APIs
<code>IBuilderConfig.set_timing_cache()</code>
<code>IGpuAllocator.reallocate()</code>
<code>INetworkDefinition.add_dequantize()</code>
<code>INetworkDefinition.add_quantize()</code>
<code>INetworkDefinition.set_weights_name()</code>
<code>IPluginRegistry.deregister_creator()</code>
<code>Refitter.get_all_weights()</code>
<code>Refitter.get_missing_weights()</code>
<code>Refitter::set_named_weights()</code>
<code>IResizeLayer.coordinate_transformation</code>
<code>IResizeLayer.nearest_rounding</code>
<code>IResizeLayer.selector_for_single_pixel</code>
<code>IScaleLayer.channel_axis</code>
<code>enum ResizeCoordinateTransformationDoc</code>
<code>enum ResizeMode</code>
<code>BuilderFlag.SPARSE_WEIGHTS</code>
<code>TacticSource.CUDNN</code>
<code>TensorFormat.DLA_HWC4</code>
<code>TensorFormat.DLA_LINEAR</code>
<code>TensorFormat.HWC16</code>

Table 13. Removed Python APIs

Removed Python APIs
Core Library
<code>class DimsCHW</code>
<code>class DimsNCHW</code>
<code>class IPlugin</code>
<code>class IPluginFactory</code>
<code>class IPluginLayer</code>
<code>class IRNNLayer</code>
<code>Builder.build_cuda_engine()</code>
<code>Builder.average_find_iterations</code>

Removed Python APIs
<code>Builder.debug_sync</code>
<code>Builder.fp16_mode</code>
<code>IBuilder.int8_mode</code>
<code>Builder.max_workspace_size</code>
<code>Builder.min_find_iterations</code>
<code>Builder.refittable</code>
<code>Builder.strict_type_constraints</code>
<code>ICudaEngine.max_workspace_size</code>
<code>IMatrixMultiplyLayer.transpose0</code>
<code>INetworkDefinition.add_matrix_multiply_deprecated()</code>
<code>INetworkDefinition.add_plugin()</code>
<code>INetworkDefinition.add_plugin_ext()</code>
<code>INetworkDefinition.add_rnn()</code>
<code>INetworkDefinition.convolution_output_dimensions_formula</code>
<code>INetworkDefinition.deconvolution_output_dimensions_formula</code>
<code>INetworkDefinition.pooling_output_dimensions_formula</code>
<code>ITensor.get_dynamic_range()</code>
<code>Dims.get_type()</code>
<code>TensorFormat.HWC8</code>
<code>TensorFormat.NCHW</code>
<code>TensorFormat.NCHW2</code>
Caffe Parser
<code>class IPluginFactory</code>
<code>class IPluginFactoryExt</code>
<code>CaffeParser.plugin_factory</code>
<code>CaffeParser.plugin_factory_ext</code>
UFF Parser
<code>class IPluginFactory</code>
<code>class IPluginFactoryExt</code>
<code>UffParser.plugin_factory</code>
<code>UffParser.plugin_factory_ext</code>

Deprecated

For our deprecation policy, refer to the [TensorRT Deprecation Policy](#) section in the *TensorRT Developer Guide*.

Table 14. Depreciated APIs

Deprecated APIs	Replaced with
<code>nvinfer1::IResizeLayer::setAlignCorners</code>	<code>nvinfer1::IResizeLayer::setCoordinateTransformation</code>
<code>nvinfer1::IResizeLayer::getAlignCorners</code>	<code>nvinfer1::IResizeLayer::setSelectorForSinglePixel</code>
	<code>nvinfer1::IResizeLayer::setNearestRounding</code>

Chapter 2. C++ API

The NVIDIA® TensorRT™ C++ API allows developers to import, calibrate, generate and deploy networks using C++. Networks can be imported directly from ONNX. They may also be created programmatically by instantiating individual layers and setting parameters and weights directly.

Within the core C++ API in `NvInfer.h`, the following APIs are included:

- ▶ [Builder API](#)
- ▶ [Execution API](#)
- ▶ [Network Definition API](#)
- ▶ [ONNX Parser API](#)
- ▶ [Plugin API](#)

To view this API, see [TensorRT C++ API](#).

For more information about the C++ API, including sample code, see [NVIDIA TensorRT Developer Guide](#).

Chapter 3. Python API

The NVIDIA® TensorRT™ Python API enables developers in Python based development environments and those looking to experiment with TensorRT to easily parse models (for example, from ONNX) and generate and run PLAN files.

To view this API, see [TensorRT Python API](#).

For more information about the Python API, including sample code, see [TensorRT Developer Guide](#).

Notice

This document is provided for information purposes only and shall not be regarded as a warranty of a certain functionality, condition, or quality of a product. NVIDIA Corporation ("NVIDIA") makes no representations or warranties, expressed or implied, as to the accuracy or completeness of the information contained in this document and assumes no responsibility for any errors contained herein. NVIDIA shall have no liability for the consequences or use of such information or for any infringement of patents or other rights of third parties that may result from its use. This document is not a commitment to develop, release, or deliver any Material (defined below), code, or functionality.

NVIDIA reserves the right to make corrections, modifications, enhancements, improvements, and any other changes to this document, at any time without notice.

Customer should obtain the latest relevant information before placing orders and should verify that such information is current and complete.

NVIDIA products are sold subject to the NVIDIA standard terms and conditions of sale supplied at the time of order acknowledgement, unless otherwise agreed in an individual sales agreement signed by authorized representatives of NVIDIA and customer ("Terms of Sale"). NVIDIA hereby expressly objects to applying any customer general terms and conditions with regards to the purchase of the NVIDIA product referenced in this document. No contractual obligations are formed either directly or indirectly by this document.

NVIDIA products are not designed, authorized, or warranted to be suitable for use in medical, military, aircraft, space, or life support equipment, nor in applications where failure or malfunction of the NVIDIA product can reasonably be expected to result in personal injury, death, or property or environmental damage. NVIDIA accepts no liability for inclusion and/or use of NVIDIA products in such equipment or applications and therefore such inclusion and/or use is at customer's own risk.

NVIDIA makes no representation or warranty that products based on this document will be suitable for any specified use. Testing of all parameters of each product is not necessarily performed by NVIDIA. It is customer's sole responsibility to evaluate and determine the applicability of any information contained in this document, ensure the product is suitable and fit for the application planned by customer, and perform the necessary testing for the application in order to avoid a default of the application or the product. Weaknesses in customer's product designs may affect the quality and reliability of the NVIDIA product and may result in additional or different conditions and/or requirements beyond those contained in this document. NVIDIA accepts no liability related to any default, damage, costs, or problem which may be based on or attributable to: (i) the use of the NVIDIA product in any manner that is contrary to this document or (ii) customer product designs.

No license, either expressed or implied, is granted under any NVIDIA patent right, copyright, or other NVIDIA intellectual property right under this document. Information published by NVIDIA regarding third-party products or services does not constitute a license from NVIDIA to use such products or services or a warranty or endorsement thereof. Use of such information may require a license from a third party under the patents or other intellectual property rights of the third party, or a license from NVIDIA under the patents or other intellectual property rights of NVIDIA.

Reproduction of information in this document is permissible only if approved in advance by NVIDIA in writing, reproduced without alteration and in full compliance with all applicable export laws and regulations, and accompanied by all associated conditions, limitations, and notices.

THIS DOCUMENT AND 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 NON-INFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE. TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL NVIDIA BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF NVIDIA HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES. Notwithstanding any damages that customer might incur for any reason whatsoever, NVIDIA's aggregate and cumulative liability towards customer for the products described herein shall be limited in accordance with the Terms of Sale for the product.

ARM

ARM, AMBA and ARM Powered are registered trademarks of ARM Limited. Cortex, MPCore and Mali are trademarks of ARM Limited. "ARM" is used to represent ARM Holdings plc; its operating company ARM Limited; and the regional subsidiaries ARM Inc.; ARM KK; ARM Korea Limited.; ARM Taiwan Limited; ARM France SAS; ARM Consulting (Shanghai) Co. Ltd.; ARM Germany GmbH; ARM Embedded Technologies Pvt. Ltd.; ARM Norway, AS and ARM Sweden AB.

HDMI

HDMI, the HDMI logo, and High-Definition Multimedia Interface are trademarks or registered trademarks of HDMI Licensing LLC.

BlackBerry/QNX

Copyright © 2020 BlackBerry Limited. All rights reserved.

Trademarks, including but not limited to BLACKBERRY, EMBLEM Design, QNX, AVIAGE, MOMENTICS, NEUTRINO and QNX CAR are the trademarks or registered trademarks of BlackBerry Limited, used under license, and the exclusive rights to such trademarks are expressly reserved.

Google

Android, Android TV, Google Play and the Google Play logo are trademarks of Google, Inc.

Trademarks

NVIDIA, the NVIDIA logo, and BlueField, CUDA, DALI, DRIVE, JetPack, Jetson AGX Xavier, Jetson Nano, Kepler, Maxwell, NGC, Nsight, Orin, Pascal, Quadro, Tegra, TensorRT, Triton, Turing and Volta are trademarks and/or registered trademarks of NVIDIA Corporation in the United States and other countries. Other company and product names may be trademarks of the respective companies with which they are associated.

Copyright

© 2017-2022 NVIDIA Corporation & affiliates. All rights reserved.

