



# NVIDIA OPTICAL FLOW SDK

## Application Note

# Table of Contents

Chapter 1. NVIDIA Optical Flow Accelerator.....	1
1.1. Introduction.....	1
1.2. NVOFA Capabilities.....	1
1.3. NVOF API.....	2
1.4. NVOFA Quality and Performance.....	3
1.5. NVOFFRUC Performance.....	5
1.6. Programming NVOFA.....	5
1.7. OpenCV Support.....	6

---

# Chapter 1. NVIDIA Optical Flow Accelerator

## 1.1. Introduction

NVIDIA® GPUs, starting with the NVIDIA Turing™ generation, contain a hardware accelerator for computing optical flow and stereo disparity between frames (referred to as NVOFA in this document), which works independently of graphics/NVIDIA CUDA® cores. With end-to-end optical flow calculation offloaded to NVOFA, the graphics/CUDA cores and the CPU are free for other operations.

Optical flow vectors are useful in various use-cases such as object detection and tracking, video frame rate up-conversion, depth estimation, stitching etc. It is also observed that using flow vectors for object detection also increases inference accuracy<sup>1</sup>. The hardware capabilities of NVOFA are exposed through APIs referred to as NVOF APIs.

## 1.2. NVOFA Capabilities

NVOFA engine can operate in two modes:

- ▶ **Optical Flow Mode:** In this mode, the engine generates flow vectors between two given frames, returning both X and Y components of the flow vectors. The Vulkan interface does not support Optical Flow mode for GPUs based on the Turing architecture.
- ▶ **Stereo Disparity Mode:** In this mode, the engine generates flow vectors in X direction only. This mode is useful in use-cases in which the Y-component of the vectors is not required, or it is known a priori that it is zero (e.g. finding disparity between the left and right images of a stereo capture). Stereo Disparity mode will be deprecated in future SDK releases. Client applications can alternatively use the X component of the flow generated in Optical Flow mode. The Vulkan interface does not support Stereo Disparity mode.

---

<sup>1</sup> Refer [http://openaccess.thecvf.com/content\\_cvpr\\_2018/papers/Yang\\_Making\\_Convolutional\\_Networks\\_CVPR\\_2018\\_paper.pdf](http://openaccess.thecvf.com/content_cvpr_2018/papers/Yang_Making_Convolutional_Networks_CVPR_2018_paper.pdf) and [http://openaccess.thecvf.com/content\\_cvpr\\_2017/papers/Carreira\\_Quo\\_Vadis\\_Action\\_CVPR\\_2017\\_paper.pdf](http://openaccess.thecvf.com/content_cvpr_2017/papers/Carreira_Quo_Vadis_Action_CVPR_2017_paper.pdf).

The hardware generates flow vectors block-wise, one vector for each block of  $4 \times 4$ ,  $2 \times 2$  and  $1 \times 1$  pixels (referred to as grid). The generated vectors can be further post-processed in software to improve accuracy; up sampled to produce dense flow map.

NVOFA hardware natively supports multiple hardware contexts with negligible context-switching penalty. As a result, subject to the hardware performance limit and available memory, an application can generate motion vectors for multiple contexts simultaneously.

The NVOFA hardware is supported for all Turing GPUs (except TU117) and above.

## 1.3. NVOF API

Capabilities of NVOFA are exposed via NVOF APIs. NVOF API includes three types of software interfaces:

- ▶ CUDA: Cross-platform API, works on Linux and Windows 10 and above.
- ▶ DirectX 11: Works on Windows 10 and above.
- ▶ DirectX 12: Works on Windows 10 20H1 and above.
- ▶ Vulkan: Cross-platform API, works on Linux and Windows 10 and above. Vulkan interface is not supported on WSL(Window subsystem for Linux) architecture.

Refer to the sample applications included in the Optical Flow SDK for more details.

[Table 1](#) and [Table 2](#) summarize the capabilities of the NVOFA hardware and the new features exposed through NVOF APIs in Optical Flow SDK 5.x respectively.

Table 1. NVOFA Hardware Capabilities

Hardware Features	Turing	Ampere	Ada
Optical flow mode	Y	Y	Y
Support for external hints	Y	Y	Y
4x4 grid size	Y	Y	Y
2x2 and 1x1 grid size	N	Y	Y
Hardware cost	N	Y	Y
Region of interest (ROI) optical flow calculation	N	Y	Y
Maximum supported resolution	4096x4096	8192x8192	8192x8192

- ▶ Y: Supported, N: Unsupported

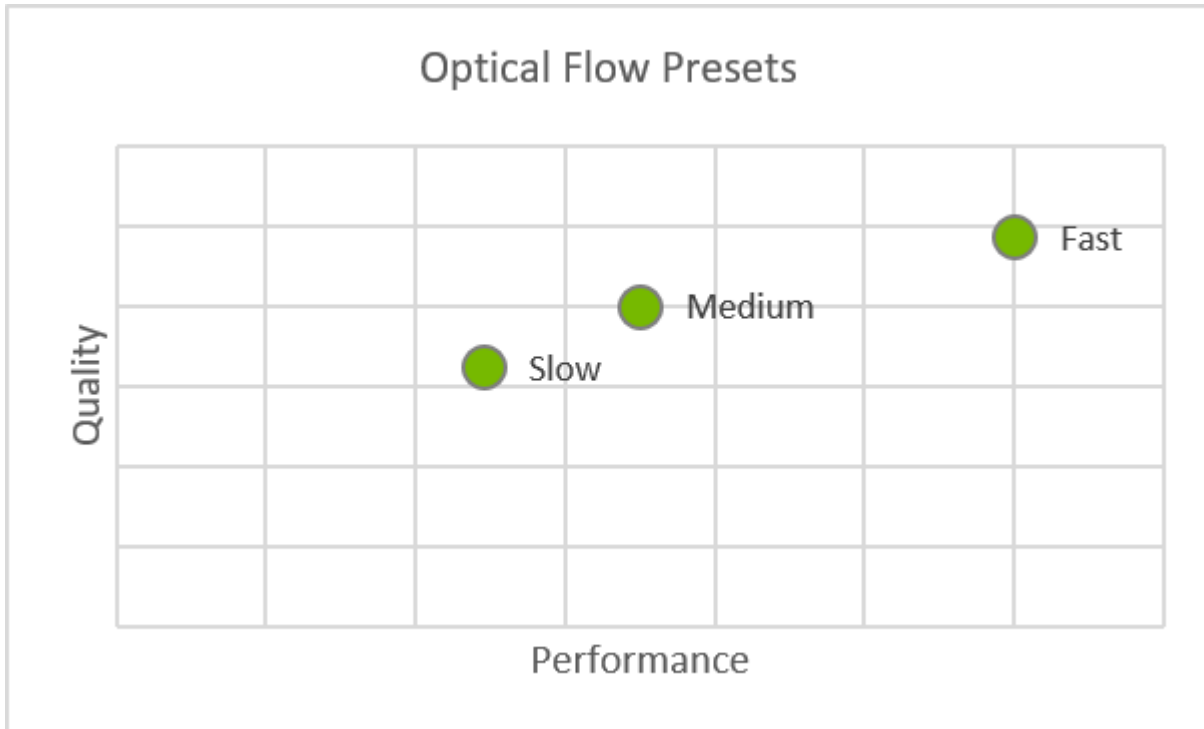
Table 2. New features exposed through NVOF SDK 5.x

Sr. No.	SDK Version	Feature	Description
1	5.0	Vulkan interface	This SDK adds support of Vulkan interface. Application can generate the flow for a pair of Vulkan resources.

## 1.4. NVOFA Quality and Performance

The NVOF API exposes multiple quality and performance levels (which are referred to as *presets*) which the user can choose based on the desired quality and performance requirement. [Figure 1](#) shows the performance/quality trade-off to be expected with the presets.

Figure 1. Performance/Quality Presets exposed in NVOF API



The NVOFA provides real time performance with small CUDA core utilization. [Table 3](#) shows *indicative*<sup>2</sup> performance in FPS with a 1080p video sequence and quality of NVOF API on KITTI 2105 which is publicly available data set. Users can trade quality vs. performance by choosing the right preset. Note that performance numbers in [Table 3](#) are measured with assumptions listed under the table. The performance varies across

<sup>2</sup> NVOFA performance depends on many factors, including but not limited to: OFAPI settings, GPU clocks, GPU type, video content type, instantaneous available memory bandwidth etc.

GPU classes (e.g. Quadro, Tesla), and scales (almost) linearly with the clock speeds for each hardware.

Table 3. Indicative quality and performance

Grid size	Preset		FI-fg			FI-all			FPS at 1080p		
			Turing	Ampere	Ada	Turing	Ampere	Ada	Turing	Ampere	Ada
4x4	SLOW	NOC	24.94	26.77	23.56	21.37	18.65	17.26	225	200	536
		OCC	27.63	29.32	26.31	31.53	29.17	27.73			
	MEDIUM	NOC	36.60	32.78	23.66	23.73	22.23	19.11	468	405	1000
		OCC	38.90	35.10	26.36	33.61	32.13	28.98			
	FAST	NOC	47.50	38.40	29.90	26.39	25.44	23.48	768	613	1296
		OCC	49.42	40.65	32.44	35.91	34.85	33.07			
2x2	SLOW	NOC	N/A	26.90	20.44	N/A	18.49	15.98	N/A	94	210
		OCC		29.43	23.18		29.08	26.00			
	MEDIUM	NOC		33.16	21.38		20.51	16.13		154	336
		OCC		35.44	24.08		30.73	26.01			
	FAST	NOC		35.44	30.81		23.09	23.13		261	711
		OCC		37.75	33.28		32.92	32.41			
1x1	SLOW	NOC		26.91	22.71		18.72	16.14		29	98
		OCC		29.45	25.44		29.16	26.42			
	MEDIUM	NOC		30.38	26.29		20.24	16.27		45	113
		OCC		32.77	28.94		30.39	26.18			
	FAST	NOC		34.55	29.88		22.65	20.93		85	222
		OCC		36.84	32.40		32.58	30.59			

- ▶ The above data is generated for Optical Flow mode using .\Samples\AppOFCuda on RTX6000, RTX3090 and RTX4090 respectively on Windows 11.
- ▶ The performance on Windows using CUDA interface with hardware scheduling disabled is typically lesser than that of Linux and Windows with hardware scheduling enabled due to a known bug inside NVIDIA display driver.
- ▶ All measurements are done by setting the video clocks as reported by nvidia-smi at 1679, 1708 and 2114MHz on RTX6000, RTX3090 and RTX4090 respectively. The performance should scale according to the actual video clocks for other GPUs. Information on nvidia-smi can be found at <https://developer.nvidia.com/nvidia-system-management-interface>.
- ▶ Resolution/Input format: 1920x1080/YUV 4:2:0
- ▶ Software: Optical flow SDK 5.0, NVIDIA display driver: 528.24 for Window OS and 525.85.05 for Linux
- ▶ FI-fg = Percentage of vectors in foreground pixels having an average EPE > 3 on KITTI 2015.

- ▶ FI-all = Percentage of vectors having an average EPE > 3 on KITTI 2015.
- ▶ Details for KITTI 2015 and the test data set can be found [http://www.cvlibs.net/datasets/kitti/eval\\_scene\\_flow.php?benchmark=flow](http://www.cvlibs.net/datasets/kitti/eval_scene_flow.php?benchmark=flow).
- ▶ NOC = Non-occluded region
- ▶ OCC = Occluded region

## 1.5. NVOFFRUC Performance

NVOFFRUC performance is measured in terms of an average time taken by NVOFFRUC library for interpolating one frame in an input sequence. It is calculated as total time taken by NVOFFRUC library to double the frame rate of input sequence divided by total number of frames in input sequence. [Table 4](#) shows *indicative*<sup>3</sup> performance as measured on publicly available dataset.

Table 4. Indicative performance

GPU Name	Performance
GeForce RTX 2080 Ti	12.01 ms
GeForce RTX 3090 Ti	9.23 ms
GeForce RTX 4090	4.41 ms

- ▶ The above data is generated with `\NvOFFRUC\NvOFFRUCSample` using NV12 CUDA array as input surface on Windows.
- ▶ All measurements are done by setting the video clocks as reported by `nvidia-smi` at 1755 MHz. The performance should scale according to the actual video clocks for other GPUs. Information on `nvidia-smi` can be found at <https://developer.nvidia.com/nvidia-system-management-interface>.
- ▶ Resolution/Input format: 1920x1080/YUV 4:2:0
- ▶ Software: Optical flow SDK 4.0 and later SDKs, NVIDIA display driver: 522.25

## 1.6. Programming NVOFA

Optical Flow SDK 5.0 is supported on R525 drivers and above. Refer to the SDK release notes for information regarding the required driver version.

Refer to the documents and the sample applications included in the SDK package for details on how to program NVOFA.

<sup>3</sup> NVOFFRUC performance depends on many factors, including but not limited to: NVOFFRUC, GPU clocks, GPU type, video content type, instantaneous available memory bandwidth etc.

## 1.7. OpenCV Support

OpenCV is one of the most popular libraries in the field of computer vision. OpenCV library contains several CPU-based and CUDA-based algorithms for computing optical flow vectors.

NVOFA can also be used with OpenCV to significantly speed up the optical flow calculation.

Note that OpenCV is an open-source project and their usage is governed by specific licenses and terms and conditions.



## 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 acknowledgment, 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.

## Trademarks

NVIDIA, the NVIDIA logo, and cuBLAS, CUDA, CUDA Toolkit, cuDNN, DALI, DIGITS, DGX, DGX-1, DGX-2, DGX Station, DLProf, GPU, Jetson, Kepler, Maxwell, NCCL, Nsight Compute, Nsight Systems, NVCAffe, NVIDIA Deep Learning SDK, NVIDIA Developer Program, NVIDIA GPU Cloud, NVLink, NVSHMEM, PerfWorks, Pascal, SDK Manager, Tegra, TensorRT, TensorRT Inference Server, Tesla, TF-TRT, Triton Inference Server, 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

© 2018-2023 NVIDIA Corporation. All rights reserved.

