



# SUPPORT MATRIX FOR TENSORRT

SWE-SWDOCTR-001-SPMT\_vTensorRT 5.1.3 | April 2019

## Support Guide



# TABLE OF CONTENTS

- Chapter 1. Features For Platforms And Software..... 1
- Chapter 2. Layers And Features..... 2
- Chapter 3. Layers And Precision.....5
- Chapter 4. Hardware And Precision.....7
- Chapter 5. Software Versions Per Platform..... 8
- Chapter 6. Supported Ops..... 9

# Chapter 1.

## FEATURES FOR PLATFORMS AND SOFTWARE

Table 1 List of supported features per platform.

	Linux x86-64	Windows x64	Linux ppc64le	Linux AArch64	QNX AArch64	Android AArch64
<a href="#">Supported CUDA versions</a>	<a href="#">9.0, 10.0, 10.1</a>	<a href="#">9.0, 10.0, 10.1</a>	<a href="#">10.1</a>	<a href="#">10.1</a>	<a href="#">10.1</a>	<a href="#">10.0</a>
<a href="#">Supported cuDNN versions</a>	<a href="#">7.5.0</a>	<a href="#">7.5.0</a>	<a href="#">7.5.0</a>	<a href="#">7.5.0</a>	<a href="#">7.5.0</a>	<a href="#">7.5.0</a>
TensorRT Python API	Yes	No	Yes	Yes	No	No
NvUffParser	Yes	Yes	Yes	Yes	Yes	Yes
NvOnnxParser	Yes	Yes	Yes	Yes	Yes	Yes



Serialized engines are not portable across platforms or TensorRT versions.

# Chapter 2.

## LAYERS AND FEATURES

Table 2 List of supported features per TensorRT layer.

Layer	Dimensions of input tensor	Dimensions of output tensor	Does the operation apply to only the innermost 3 dimensions?	Supports broadcast <sup>1</sup>	Supports broadcast across batch <sup>2</sup>
<a href="#">Activation</a>	0-7 dimensions	0-7 dimensions	No	No	No
<a href="#">Concatenation</a>	1-7 dimensions	1-7 dimensions	No	No	No
<a href="#">Constant</a>	0-7 dimensions	0-7 dimensions	No	No	Always
<a href="#">Convolution</a>	3 or more dimensions	3 or more dimensions	Yes	No	No
<a href="#">Deconvolution</a>	3 or more dimensions	3 or more dimensions	Yes	No	No
<a href="#">ElementWise</a>	0-7 dimensions	0-7 dimensions	No	Yes	Yes
<a href="#">FullyConnected</a>	3 or more dimensions	3 or more dimensions	Yes	No	No
<a href="#">Gather</a>	<ul style="list-style-type: none"> <li>▶ Input1: 1-7 dimensions</li> <li>▶ Input2: 0-7 dimensions</li> </ul>	0-7 dimensions	No	No	Yes
<a href="#">Identity</a>	0-7 dimensions	0-7 dimensions	No	No	No

Layer	Dimensions of input tensor	Dimensions of output tensor	Does the operation apply to only the innermost 3 dimensions?	Supports broadcast <sup>1</sup>	Supports broadcast across batch <sup>2</sup>
<a href="#">IPluginV2</a>	User defined	User defined	User defined	User defined	User defined
<a href="#">LRN</a>	3 or more dimensions	3 or more dimensions	Yes	No	No
<a href="#">MatrixMultiply</a>	2 or more dimensions	2 or more dimensions	No	Yes	Yes
<a href="#">Padding</a>	3 or more dimensions	3 or more dimensions	Yes	No	No
<a href="#">Plugin</a>	User defined	User defined	User defined	User defined	User defined
<a href="#">Pooling</a>	3 or more dimensions	3 or more dimensions	Yes	Yes	Yes
<a href="#">RaggedSoftMax</a>	<ul style="list-style-type: none"> <li>▶ Input: 2 dimensions</li> <li>▶ Bounds: 2 dimensions</li> </ul>	2 or more dimensions	No	No	Yes
<a href="#">Reduce</a>	1-7 dimensions	0-7 dimensions	No	No	No
<a href="#">RNN</a>	3 dimensions	3 dimensions	No	No	No
<a href="#">RNNv2</a>	<ul style="list-style-type: none"> <li>▶ Data/Hidden/Cell: 2 or more dimensions</li> <li>▶ SeqLen: 0 or more dimensions</li> </ul>	Data/Hidden/Cell: 2 or more dimensions	No	No	No
<a href="#">Scale</a>	3 or more dimensions	3 or more dimensions	Yes	No	No
<a href="#">Shuffle</a>	0-7 dimensions	0-7 dimensions	No	No	No
<a href="#">Slice</a>	1-7 dimensions	1-7 dimensions	No	No	Yes
<a href="#">SoftMax</a>	1-7 dimensions	1-7 dimensions	No	No	Yes

Layer	Dimensions of input tensor	Dimensions of output tensor	Does the operation apply to only the innermost 3 dimensions?	Supports broadcast <sup>1</sup>	Supports broadcast across batch <sup>2</sup>
<a href="#">TopK</a>	1-7 dimensions	<ul style="list-style-type: none"> <li>▶ Output1: 1-7 dimensions</li> <li>▶ Output2: 1-7 dimensions</li> </ul>	Yes	No	Yes
<a href="#">Unary</a>	0-7 dimensions	0-7 dimensions	No	No	No

For more information about each of the TensorRT layers, see [TensorRT Layers](#).

# Chapter 3.

## LAYERS AND PRECISION

The following table lists the TensorRT layers and the precision modes that each layer supports. It also lists the ability of the layer to run on Deep Learning Accelerator (DLA). For more information about additional constraints, see [DLA Supported Layers](#).

For more information about each of the TensorRT layers, see [TensorRT Layers](#). To view a list of the specific attributes that are supported by each layer, refer to the [TensorRT API](#) documentation.

Table 3 List of supported precision mode per TensorRT layer.

Layer	FP32	FP16	INT8	DLA <sup>3</sup>
<a href="#">Activation</a>	Yes	Yes	Yes	Yes
<a href="#">Concatenation</a>	Yes	Yes	Yes	Yes
<a href="#">Constant</a>	Yes	Yes	Yes	No
<a href="#">Convolution</a>	Yes	Yes	Yes	Yes
<a href="#">Deconvolution</a>	Yes	Yes	Yes	Yes
<a href="#">ElementWise</a>	Yes	Yes	No	Yes
<a href="#">FullyConnected</a>	Yes	Yes	Yes	Yes
<a href="#">Gather</a>	Yes	Yes	No	No
<a href="#">Identity</a>	Yes	Yes	Yes	No
<a href="#">IPluginV2</a>	Yes	Yes	No	No
<a href="#">LRN</a>	Yes	Yes	Yes	Yes
<a href="#">MatrixMultiply</a>	Yes	Yes	No	No
<a href="#">Padding</a>	Yes	Yes	Yes	No

Layer	FP32	FP16	INT8	DLA <sup>3</sup>
<a href="#">Plugin</a>	Yes	Yes	No	No
<a href="#">Pooling</a>	Yes	Yes	Yes	Yes
<a href="#">RaggedSoftMax</a>	Yes	No	No	No
<a href="#">Reduce</a>	Yes	Yes	No	No
<a href="#">RNN</a>	Yes	Yes	No	No
<a href="#">RNNv2</a>	Yes	Yes	No	No
<a href="#">Scale</a>	Yes	Yes	Yes	Yes
<a href="#">Shuffle</a>	Yes	Yes	Yes	No
<a href="#">Slice</a>	Yes	Yes	No <sup>4</sup>	No
<a href="#">SoftMax</a>	Yes	Yes	No	No
<a href="#">TopK</a>	Yes	Yes	No	No
<a href="#">Unary</a>	Yes	Yes	No	No

<sup>4</sup> Partial support. Yes for unstrided Slice and No for strided.



# Chapter 4.

## HARDWARE AND PRECISION

The following table lists NVIDIA hardware and which precision modes each hardware supports. It also lists availability of Deep Learning Accelerator (DLA) on these hardware. TensorRT supports all NVIDIA hardware with capability SM 3.0 or higher.

Table 4 List of supported precision mode per hardware.

<a href="#">CUDA Compute Capability</a>	Example Device	FP32	FP16	INT8	FP16 Tensor Cores	INT8 Tensor Cores	DLA
7.5	Tesla T4	Yes	Yes	Yes	Yes	Yes	No
7.2	Jetson AGX Xavier	Yes	Yes	Yes	Yes	Yes	Yes
7.0	Tesla V100	Yes	Yes	Yes	Yes	No	No
6.2	Jetson TX2	Yes	Yes	No	No	No	No
6.1	Tesla P4	Yes	No	Yes	No	No	No
6.0	Tesla P100	Yes	Yes	No	No	No	No
5.3	Jetson TX1	Yes	Yes	No	No	No	No
5.2	Tesla M4	Yes	No	No	No	No	No
5.0	Quadro K2200	Yes	No	No	No	No	No
3.7	Tesla K80	Yes	No	No	No	No	No
3.5	Tesla K40	Yes	No	No	No	No	No
3.0	Tesla K10	Yes	No	No	No	No	No

# Chapter 5.

## SOFTWARE VERSIONS PER PLATFORM

Table 5 List of supported platforms per software version.

	Compiler version	Python version
Ubuntu 14.04 x86-64	<a href="#">gcc 4.8.4</a>	2.7, 3.4
Ubuntu 16.04 x86-64	<a href="#">gcc 5.4.0</a>	2.7, 3.5
Ubuntu 18.04 x86-64	<a href="#">gcc 7.3.0</a>	2.7, 3.6
CentOS 7.5 x86-64	<a href="#">gcc 4.8.5</a>	2.7, 3.6
Windows 10 x64	<b>CUDA 10.0, 10.1</b> <a href="#">MSVC 2017u5</a> <b>CUDA 9.0</b> <a href="#">MSVC 2017u3</a>	
Ubuntu 18.04 ppc64le	<a href="#">gcc 7.3.0</a>	2.7, 3.6
CentOS 7.5 ppc64le	<a href="#">gcc 4.8.5</a>	2.7, 3.6
Ubuntu 18.04 AArch64	<a href="#">gcc 7.3.1</a>	2.7, 3.6
QNX AArch64	<a href="#">gcc 5.4.0</a>	
Android AArch64	<a href="#">Clang 5.0.300080</a>	

# Chapter 6.

## SUPPORTED OPS

The following lists describe the operations that are supported in a Caffe or TensorFlow framework and in the ONNX TensorRT parser:

### Caffe

These are the operations that are supported in a Caffe framework:

- ▶ **BatchNormalization**
- ▶ **BNLL**
- ▶ **Clip**<sup>5</sup>
- ▶ **Concatenation**
- ▶ **Convolution**
- ▶ **Crop**
- ▶ **Deconvolution**
- ▶ **Dropout**
- ▶ **ElementWise**
- ▶ **ELU**
- ▶ **InnerProduct**
- ▶ **Input**
- ▶ **LeakyReLU**
- ▶ **LRN**
- ▶ **Permute**
- ▶ **Pooling**
- ▶ **Power**
- ▶ **Reduction**
- ▶ **ReLU, TanH, and Sigmoid**
- ▶ **Reshape**

---

<sup>5</sup> When using the `Clip` operation, Caffe users must serialize their layers using `ditcaffe.pb.h` instead of `caffe.pb.h` in order to import the layer into TensorRT.

- ▶ **SoftMax**
- ▶ **Scale**

## TensorFlow

These are the operations that are supported in a TensorFlow framework:

- ▶ **Add, Sub, Mul, Div, Minimum and Maximum**
- ▶ **ArgMax**
- ▶ **ArgMin**
- ▶ **AvgPool**
- ▶ **BiasAdd**
- ▶ **Clip**
- ▶ **ConcatV2**
- ▶ **Const**
- ▶ **Conv2D**
- ▶ **ConvTranspose2D**
- ▶ **DepthwiseConv2dNative**
- ▶ **Elu**
- ▶ **ExpandDims**
- ▶ **FusedBatchNorm**
- ▶ **Identity**
- ▶ **LeakyReLU**
- ▶ **MaxPool**
- ▶ **Mean**
- ▶ **Negative, Abs, Sqrt, Recip, Rsqrt, Pow, Exp and Log**
- ▶ **Pad** is supported if followed by one of these TensorFlow layers: **Conv2D**, **DepthwiseConv2dNative**, **MaxPool**, and **AvgPool**.
- ▶ **Placeholder**
- ▶ **ReLU, TanH, and Sigmoid**
- ▶ **Relu6**
- ▶ **Reshape**
- ▶ **Sin, Cos, Tan, Asin, Acos, Atan, Sinh, Cosh, Asinh, Acosh, Atanh, Ceil and Floor**
- ▶ **Selu**
- ▶ **Slice**
- ▶ **SoftMax**



If the input to a TensorFlow `softMax` op is not `NHWC`, TensorFlow will automatically insert a transpose layer with a non-constant permutation, causing

the UFF converter to fail. It is therefore advisable to manually transpose `SoftMax` inputs to `NHWC` using a constant permutation.

- ▶ `Softplus`
- ▶ `Softsign`
- ▶ `Transpose`

## ONNX

Since the ONNX parser is an open source project, the most up-to-date information regarding the supported operations can be found in [GitHub: ONNX TensorRT](#).

These are the operations that are supported in the ONNX framework:

- ▶ `Abs`
- ▶ `Add`
- ▶ `ArgMax`
- ▶ `ArgMin`
- ▶ `AveragePool`
- ▶ `BatchNormalization`
- ▶ `Cast`
- ▶ `Ceil`
- ▶ `Clip`
- ▶ `Concat`
- ▶ `Constant`
- ▶ `Conv`
- ▶ `ConvTranspose`
- ▶ `DepthToSpace`
- ▶ `Div`
- ▶ `Dropout`
- ▶ `Elu`
- ▶ `Exp`
- ▶ `Flatten`
- ▶ `Floor`
- ▶ `Gather`
- ▶ `Gemm`
- ▶ `GlobalAveragePool`
- ▶ `GlobalMaxPool`
- ▶ `HardSigmoid`
- ▶ `Identity`
- ▶ `ImageScaler`
- ▶ `InstanceNormalization`
- ▶ `LRN`

- ▶ LeakyRelU
- ▶ Log
- ▶ LogSoftmax
- ▶ MatMul
- ▶ Max
- ▶ MaxPool
- ▶ Mean
- ▶ Min
- ▶ Mul
- ▶ Neg
- ▶ Pad
- ▶ ParametricSoftplus
- ▶ Pow
- ▶ Reciprocal
- ▶ ReduceL1
- ▶ ReduceL2
- ▶ ReduceLogSum
- ▶ ReduceLogSumExp
- ▶ ReduceMax
- ▶ ReduceMean
- ▶ ReduceMin
- ▶ ReduceProd
- ▶ ReduceSum
- ▶ ReduceSumSquare
- ▶ Relu
- ▶ Reshape
- ▶ ScaledTanh
- ▶ Selu
- ▶ Shape
- ▶ Sigmoid
- ▶ Sin, Cos, Tan, Asin, Acos, Atan, Sinh, Cosh, Asinh, Acosh, and Atanh
- ▶ Size
- ▶ Slice
- ▶ Softmax
- ▶ Softplus
- ▶ Softsign
- ▶ SpaceToDepth
- ▶ Split
- ▶ Squeeze

- ▶ **Sub**
- ▶ **Sum**
- ▶ **Tanh**
- ▶ **ThresholdedRelu**
- ▶ **TopK**
- ▶ **Transpose**
- ▶ **Unsqueeze**
- ▶ **Upsample**

## Notice

THE INFORMATION IN THIS GUIDE AND ALL OTHER INFORMATION CONTAINED IN NVIDIA DOCUMENTATION REFERENCED IN THIS GUIDE IS PROVIDED “AS IS.” NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO THE INFORMATION FOR THE PRODUCT, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE. Notwithstanding any damages that customer might incur for any reason whatsoever, NVIDIA’s aggregate and cumulative liability towards customer for the product described in this guide shall be limited in accordance with the NVIDIA terms and conditions of sale for the product.

THE NVIDIA PRODUCT DESCRIBED IN THIS GUIDE IS NOT FAULT TOLERANT AND IS NOT DESIGNED, MANUFACTURED OR INTENDED FOR USE IN CONNECTION WITH THE DESIGN, CONSTRUCTION, MAINTENANCE, AND/OR OPERATION OF ANY SYSTEM WHERE THE USE OR A FAILURE OF SUCH SYSTEM COULD RESULT IN A SITUATION THAT THREATENS THE SAFETY OF HUMAN LIFE OR SEVERE PHYSICAL HARM OR PROPERTY DAMAGE (INCLUDING, FOR EXAMPLE, USE IN CONNECTION WITH ANY NUCLEAR, AVIONICS, LIFE SUPPORT OR OTHER LIFE CRITICAL APPLICATION). NVIDIA EXPRESSLY DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY OF FITNESS FOR SUCH HIGH RISK USES. NVIDIA SHALL NOT BE LIABLE TO CUSTOMER OR ANY THIRD PARTY, IN WHOLE OR IN PART, FOR ANY CLAIMS OR DAMAGES ARISING FROM SUCH HIGH RISK USES.

NVIDIA makes no representation or warranty that the product described in this guide will be suitable for any specified use without further testing or modification. Testing of all parameters of each product is not necessarily performed by NVIDIA. It is customer’s sole responsibility to ensure the product is suitable and fit for the application planned by customer and to do 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 guide. NVIDIA does not accept any 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 guide, or (ii) customer product designs.

Other than the right for customer to use the information in this guide with the product, no other license, either expressed or implied, is hereby granted by NVIDIA under this guide. Reproduction of information in this guide is permissible only if reproduction is approved by NVIDIA in writing, is reproduced without alteration, and is accompanied by all associated conditions, limitations, and notices.

## Trademarks

NVIDIA, the NVIDIA logo, and cuBLAS, CUDA, cuDNN, cuFFT, cuSPARSE, DALI, DIGITS, DGX, DGX-1, Jetson, Kepler, NVIDIA Maxwell, NCCL, NVLink, Pascal, Tegra, TensorRT, and Tesla 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

© 2019 NVIDIA Corporation. All rights reserved.