



# NVIDIA TensorRT

API Migration Guide | NVIDIA Docs

# Table of Contents

Chapter 1. Python.....	1
1.1. Python API Changes.....	1
1.2. Added Python APIs.....	3
1.3. Removed Python APIs.....	4
Chapter 2. C++.....	7
2.1. C++ API Changes.....	7
2.2. 64-Bit Dimension Changes.....	8
2.3. Added C++ APIs.....	8
2.4. Removed C++ APIs.....	10
2.5. Removed C++ Plugins.....	14
2.6. Removed Safety C++ APIs.....	15
Chapter 3. trtexec.....	16
3.1. trtexec Flag Changes.....	16
3.2. Removed trtexec Flags.....	16
3.3. Deprecated trtexec Flags.....	16

# Chapter 1. Python

## 1.1. Python API Changes

Table 1. Allocating Buffers and Using a Name-Based Engine API

TensorRT 8.x	TensorRT 10.0
<pre>def allocate_buffers(self, engine):     """     Allocates all buffers required for     an engine, i.e. host/device inputs/     outputs.     """     inputs = []     outputs = []     bindings = []     stream = cuda.Stream()      # binding is the name of input/     output     for binding in engine:         size =     trt.volume(engine.get_binding_shape(binding)     * engine.max_batch_size         dtype =     trt.nptype(engine.get_binding_dtype(binding)          # Allocate host and device         buffers         host_mem =     cuda.pagelocked_empty(size, dtype)     # page-locked memory buffer (won't     swapped to disk)         device_mem =     cuda.mem_alloc(host_mem.nbytes)          # Append the device buffer         address to device bindings.         # When cast to int, it's a         linear index into the context's memory         (like memory address).         bindings.append(int(device_mem))</pre>	<pre>def allocate_buffers(self, engine):     """     Allocates all buffers required for     an engine, i.e. host/device inputs/     outputs.     """     inputs = []     outputs = []     bindings = []     stream = cuda.Stream()      for i in     range(engine.num_io_tensors):         tensor_name =     engine.get_tensor_name(i)         size =     trt.volume(engine.get_tensor_shape(tensor_name))         dtype =     trt.nptype(engine.get_tensor_dtype(tensor_name))          # Allocate host and device         buffers         host_mem =     cuda.pagelocked_empty(size, dtype)     # page-locked memory buffer (won't     swapped to disk)         device_mem =     cuda.mem_alloc(host_mem.nbytes)          # Append the device buffer         address to device bindings.         # When cast to int, it's a         linear index into the context's memory         (like memory address).         bindings.append(int(device_mem))          # Append to the appropriate         input/output list.</pre>

TensorRT 8.x	TensorRT 10.0
<pre># Append to the appropriate input/output list. if engine.binding_is_input(binding):  inputs.append(self.HostDeviceMem(host_mem, device_mem)) else:  outputs.append(self.HostDeviceMem(host_mem, device_mem))  return inputs, outputs, bindings, stream</pre>	<pre>if engine.get_tensor_mode(tensor_name) == trt.TensorIOMode.INPUT:  inputs.append(self.HostDeviceMem(host_mem, device_mem)) else:  outputs.append(self.HostDeviceMem(host_mem, device_mem))  return inputs, outputs, bindings, stream</pre>

Table 2. Transition from enqueueV2 to enqueueV3 for Python

TensorRT 8.x	TensorRT 10.0
<pre># Allocate device memory for inputs. d_inputs = [cuda.mem_alloc(input_nbytes) for binding in range(input_num)]  # Allocate device memory for outputs. h_output = cuda.pagelocked_empty(output_nbytes, dtype=np.float32) d_output = cuda.mem_alloc(h_output.nbytes)  # Transfer data from host to device. cuda.memcpy_htod_async(d_inputs[0], input_a, stream) cuda.memcpy_htod_async(d_inputs[1], input_b, stream) cuda.memcpy_htod_async(d_inputs[2], input_c, stream)  # Run inference context.execute_async_v2(bindings=[int(d_input) for d_inp in d_inputs] + [int(d_output)], stream_handle=stream.handle)  # Synchronize the stream stream.synchronize()</pre>	<pre># Allocate device memory for inputs. d_inputs = [cuda.mem_alloc(input_nbytes) for binding in range(input_num)]  # Allocate device memory for outputs. h_output = cuda.pagelocked_empty(output_nbytes, dtype=np.float32) d_output = cuda.mem_alloc(h_output.nbytes)  # Transfer data from host to device. cuda.memcpy_htod_async(d_inputs[0], input_a, stream) cuda.memcpy_htod_async(d_inputs[1], input_b, stream) cuda.memcpy_htod_async(d_inputs[2], input_c, stream)  # Setup tensor address bindings = [int(d_inputs[i]) for i in range(3)] + [int(d_output)]  for i in range(engine.num_io_tensors):  context.set_tensor_address(engine.get_tensor_name(i), bindings[i])  # Run inference context.execute_async_v3(stream_handle=stream.handle)  # Synchronize the stream stream.synchronize()</pre>

Table 3. Engine Building, use only build\_serialized\_network

TensorRT 8.x	TensorRT 10.0
<pre>engine_bytes = None try:</pre>	<pre>engine_bytes = self.builder.build_serialized_network(self.network, self.config)</pre>

TensorRT 8.x	TensorRT 10.0
<pre> engine_bytes = self.builder.build_serialized_network(self.network, self.config) except AttributeError: engine = self.builder.build_engine(self.network, self.config) engine_bytes = engine.serialize() del engine assert engine_bytes </pre>	<pre> if engine_bytes is None: log.error("Failed to create engine") sys.exit(1) </pre>

## 1.2. Added Python APIs

### Types

- ▶ APILanguage
- ▶ ExecutionContextAllocationStrategy
- ▶ IGpuAsyncAllocator
- ▶ InterfaceInfo
- ▶ IPluginResource
- ▶ IPluginV3
- ▶ IStreamReader
- ▶ IVersionedInterface

### Methods and Properties

- ▶ ICudaEngine.is\_debug\_tensor()
- ▶ ICudaEngine.minimum\_weight\_streaming\_budget
- ▶ ICudaEngine.streamable\_weights\_size
- ▶ ICudaEngine.weight\_streaming\_budget
- ▶ IExecutionContext.get\_debug\_listener()
- ▶ IExecutionContext.get\_debug\_state()
- ▶ IExecutionContext.set\_all\_tensors\_debug\_state()
- ▶ IExecutionContext.set\_debug\_listener()
- ▶ IExecutionContext.set\_tensor\_debug\_state()
- ▶ IExecutionContext.update\_device\_memory\_size\_for\_shapes()
- ▶ IGpuAllocator.allocate\_async()
- ▶ IGpuAllocator.deallocate\_async()
- ▶ INetworkDefinition.add\_plugin\_v3()
- ▶ INetworkDefinition.is\_debug\_tensor()

- ▶ `INetworkDefinition.mark_debug()`
- ▶ `INetworkDefinition.unmark_debug()`
- ▶ `IPluginRegistry.acquire_plugin_resource()`
- ▶ `IPluginRegistry.all_creators`
- ▶ `IPluginRegistry.deregister_creator()`
- ▶ `IPluginRegistry.get_creator()`
- ▶ `IPluginRegistry.register_creator()`
- ▶ `IPluginRegistry.release_plugin_resource()`

## 1.3. Removed Python APIs

Table 4. Removed Python APIs and their Suggested Superseded API

Python API	Superseded API
<code>BuilderFlag.ENABLE_TACTIC_HEURISTIC</code>	Builder optimization level 2
<code>BuilderFlag.STRICT_TYPES</code>	Use all three flags: <ol style="list-style-type: none"> <li>1. <code>BuilderFlag.DIRECT_IO</code></li> <li>2. <code>BuilderFlag.PREFER_PRECISION_CONSTRAINTS</code></li> <li>3. <code>BuilderFlag.REJECT_EMPTY_ALGORITHMS</code></li> </ol>
<ol style="list-style-type: none"> <li>1. <code>EngineCapability.DEFAULT</code></li> <li>2. <code>EngineCapability.kSAFE_DLA</code></li> <li>3. <code>EngineCapability.SAFE_GPU</code></li> </ol>	<ol style="list-style-type: none"> <li>1. <code>EngineCapability.STANDARD</code></li> <li>2. <code>EngineCapability.DLA_STANDALONE</code></li> <li>3. <code>EngineCapability.SAFETY</code></li> </ol>
<code>IAlgorithmIOInfo.tensor_format</code>	The strides, data type, and vectorization information is sufficient to uniquely identify tensor formats.
<code>IBuilder.max_batch_size</code>	Implicit batch is no longer supported.
<code>IBuilderConfig.max_workspace_size</code>	<ol style="list-style-type: none"> <li>1. <code>IBuilderConfig.set_memory_pool_limit()</code> with <code>MemoryPoolType.WORKSPACE</code></li> <li>2. <code>IBuilderConfig.get_memory_pool_limit()</code> with <code>MemoryPoolType.WORKSPACE</code></li> </ol>
<code>IBuilderConfig.min_timing_iterations</code>	<code>IBuilderConfig.avg_timing_iterations</code>
<ol style="list-style-type: none"> <li>1. <code>ICudaEngine.binding_is_input()</code></li> <li>2. <code>ICudaEngine.get_binding_bytes_per_component()</code></li> <li>3. <code>ICudaEngine.get_binding_components_per_element()</code></li> </ol>	<ol style="list-style-type: none"> <li>1. <code>ICudaEngine.get_tensor_mode()</code></li> <li>2. <code>ICudaEngine.get_tensor_bytes_per_component()</code></li> <li>3. <code>ICudaEngine.get_tensor_components_per_element()</code></li> </ol>

Python API	Superseded API
<ul style="list-style-type: none"> <li>4. ICudaEngine.get_binding_dtype()</li> <li>5. ICudaEngine.get_binding_format()</li> <li>6. ICudaEngine.get_binding_format_desc()</li> <li>7. ICudaEngine.get_binding_index()</li> <li>8. ICudaEngine.get_binding_name()</li> <li>9. ICudaEngine.get_binding_shape()</li> <li>10. ICudaEngine.get_binding_vectorized_dimensions()</li> <li>11. ICudaEngine.get_location()</li> <li>12. ICudaEngine.get_profile_shape()</li> <li>13. ICudaEngine.get_profile_shape_input()</li> <li>14. ICudaEngine.has_implicit_batch_dimension()</li> <li>15. ICudaEngine.is_execution_binding()</li> <li>16. ICudaEngine.is_shape_binding()</li> <li>17. ICudaEngine.max_batch_size()</li> <li>18. ICudaEngine.num_bindings()</li> </ul>	<ul style="list-style-type: none"> <li>4. ICudaEngine.get_tensor_dtype()</li> <li>5. ICudaEngine.get_tensor_format()</li> <li>6. ICudaEngine.get_tensor_format_desc()</li> <li>7. No name-based equivalent replacement</li> <li>8. No name-based equivalent replacement</li> <li>9. ICudaEngine.get_tensor_shape()</li> <li>10. ICudaEngine.get_tensor_vectorized_dimensions()</li> <li>11. ITensor.location</li> <li>12. ICudaEngine.get_tensor_profile_shape()</li> <li>13. ICudaEngine.get_tensor_profile_values()</li> <li>14. Implicit batch is no longer supported</li> <li>15. No name-based equivalent replacement</li> <li>16. ICudaEngine.is_shape_inference_io()</li> <li>17. Implicit batch is no longer supported</li> <li>18. ICudaEngine.num_io_tensors()</li> </ul>
<ul style="list-style-type: none"> <li>1. IExecutionContext.get_binding_shape()</li> <li>2. IExecutionContext.get_strides()</li> <li>3. IExecutionContext.set_binding_shape()</li> </ul>	<ul style="list-style-type: none"> <li>1. IExecutionContext.get_tensor_shape()</li> <li>2. IExecutionContext.get_tensor_strides()</li> <li>3. IExecutionContext.set_input_shape()</li> </ul>
IFullyConnectedLayer	IMatrixMultiplyLayer
<ul style="list-style-type: none"> <li>1. INetworkDefinition.add_convolution()</li> <li>2. INetworkDefinition.add_deconvolution()</li> <li>3. INetworkDefinition.add_fully_connected()</li> <li>4. INetworkDefinition.add_padding()</li> <li>5. INetworkDefinition.add_pooling()</li> <li>6. INetworkDefinition.add_rnn_v2()</li> <li>7. INetworkDefinition.has_explicit_precision</li> <li>8. INetworkDefinition.has_implicit_batch_dimension</li> </ul>	<ul style="list-style-type: none"> <li>1. INetworkDefinition.add_convolution_nd()</li> <li>2. INetworkDefinition.add_deconvolution_nd()</li> <li>3. INetworkDefinition.add_matrix_multiply()</li> <li>4. INetworkDefinition.add_padding_nd()</li> <li>5. INetworkDefinition.add_pooling_nd()</li> <li>6. INetworkDefinition.add_loop()</li> <li>7. Explicit precision support is removed in 10.0</li> <li>8. Implicit batch is no longer supported</li> </ul>
IRNNv2Layer	ILoop
<ul style="list-style-type: none"> <li>1. NetworkDefinitionCreationFlag.EXPLICIT_BATCH</li> <li>2. NetworkDefinitionCreationFlag.EXPLICIT_PRECISION</li> </ul>	<p>Support is removed in 10.0</p>
<ul style="list-style-type: none"> <li>1. PaddingMode.CAFFE_ROUND_DOWN</li> </ul>	<p>Caffe is not supported since 9.0</p>

Python API	Superseded API
2. <code>PaddingMode.CAFFE_ROUND_UP</code>	
1. <code>PreviewFeature.DISABLE_EXTERNAL_TACTICS_3005</code> 2. <code>PreviewFeature.FASTER_DYNAMIC_SHAPES_0805</code>	1. External tactics are always disabled for core code 2. This flag is on by default
1. <code>ProfilingVerbosity.DEFAULT</code> 2. <code>ProfilingVerbosity.VERBOSE</code>	1. <code>ProfilingVerbosity.LAYER_NAMES_ONLY</code> 2. <code>ProfilingVerbosity.DETAILED</code>
ResizeMode	Use <code>InterpolationMode</code> , alias is removed
SampleMode.DEFAULT	SampleMode.STRICT_BOUNDS
SliceMode	Use <code>SampleMode</code> , alias is removed



# Chapter 2. C++

## 2.1. C++ API Changes

Table 5. Transition from enqueueV2 to enqueueV3 for C++

TensorRT 8.x	TensorRT 10.0
<pre>// Create RAII buffer manager object. samplesCommon::BufferManager   buffers(mEngine);  auto context =   SampleUniquePtr&lt;nvinfer1::IExecutionContext&gt;::createExecutionContext(); if (!context) {   return false; }  // Pick a random digit to try to infer. srand(time(NULL)); int32_t const digit = rand() % 10;  // Read the input data into the managed   buffers. // There should be just 1 input tensor. ASSERT(mParams.inputTensorNames.size()   == 1);  if (!processInput(buffers,   mParams.inputTensorNames[0], digit)) {   return false; } // Create CUDA stream for the execution   of this inference. cudaStream_t stream; CHECK(cudaStreamCreate(&amp;stream));</pre>	<pre>// Create RAII buffer manager object. samplesCommon::BufferManager   buffers(mEngine);  auto context =   SampleUniquePtr&lt;nvinfer1::IExecutionContext&gt;(mEngine- &gt;createExecutionContext()); if (!context) {   return false; }  for (int32_t i = 0, e = mEngine- &gt;getNbIOTensors(); i &lt; e; i++) {   auto const name = mEngine- &gt;getIOTensorName(i);   context-&gt;setTensorAddress(name,   buffers.getDeviceBuffer(name)); }  // Pick a random digit to try to infer. srand(time(NULL)); int32_t const digit = rand() % 10;  // Read the input data into the managed   buffers. // There should be just 1 input tensor. ASSERT(mParams.inputTensorNames.size()   == 1);  if (!processInput(buffers,   mParams.inputTensorNames[0], digit)) {   return false; } // Create CUDA stream for the execution   of this inference.</pre>

TensorRT 8.x	TensorRT 10.0
<pre> // Asynchronously copy data from host input buffers to device input buffers buffers.copyInputToDeviceAsync(stream);  // Asynchronously enqueue the inference work  if (!context- &gt;enqueueV2(buffers.getDeviceBindings().data; stream, nullptr)) {     return false; } // Asynchronously copy data from device output buffers to host output buffers. buffers.copyOutputToHostAsync(stream);  // Wait for the work in the stream to complete. CHECK(cudaStreamSynchronize(stream));  // Release stream. CHECK(cudaStreamDestroy(stream)); </pre>	<pre> cudaStream_t stream; CHECK(cudaStreamCreate(&amp;stream));  // Asynchronously copy data from host input buffers to device input buffers buffers.copyInputToDeviceAsync(stream);  // Asynchronously enqueue the inference workif (!context-&gt;enqueueV3(stream)) {     return false; } // Asynchronously copy data from device output buffers to host output buffers. buffers.copyOutputToHostAsync(stream);  // Wait for the work in the stream to complete. CHECK(cudaStreamSynchronize(stream));  // Release stream. CHECK(cudaStreamDestroy(stream)); </pre>

## 2.2. 64-Bit Dimension Changes

The dimensions held by Dims changed from `int32_t` to `int64_t`. However, in TensorRT 10.0, TensorRT will generally reject networks that use dimensions exceeding the range of `int32_t`. The tensor type returned by `IShapeLayer` is now `DataType::kINT64`. Use `ICastLayer` to cast the result to the tensor of type `DataType::kINT32` if 32-bit dimensions are required.

Inspect code that bitwise copies to and from Dims to ensure it is correct for `int64_t` dimensions.

## 2.3. Added C++ APIs

### Enums

- ▶ `ActivationType::kGELU_ERF`
- ▶ `ActivationType::kGELU_TANH`
- ▶ `BuilderFlag::kREFIT_IDENTICAL`
- ▶ `BuilderFlag::kSTRIP_PLAN`
- ▶ `BuilderFlag::kWEIGHT_STREAMING`
- ▶ `Datatype::kINT4`
- ▶ `LayerType::kPLUGIN_V3`

## Types

- ▶ `APILanguage`
- ▶ `Dims64`
- ▶ `ExecutionContextAllocationStrategy`
- ▶ `IGpuAsyncAllocator`
- ▶ `InterfaceInfo`
- ▶ `IPluginResource`
- ▶ `IPluginV3`
- ▶ `IStreamReader`
- ▶ `IVersionedInterface`


## Methods and Properties

- ▶ `getInferLibBuildVersion`
- ▶ `getInferLibMajorVersion`
- ▶ `getInferLibMinorVersion`
- ▶ `getInferLibPatchVersion`
- ▶ `ICudaEngine::createRefitter`
- ▶ `IcudaEngine::getMinimumWeightStreamingBudget`
- ▶ `IcudaEngine::getStreamableWeightsSize`
- ▶ `ICudaEngine::getWeightStreamingBudget`
- ▶ `IcudaEngine::isDebugTensor`
- ▶ `ICudaEngine::setWeightStreamingBudget`
- ▶ `IExecutionContext::getDebugListener`
- ▶ `IExecutionContext::getTensorDebugState`
- ▶ `IExecutionContext::setAllTensorsDebugState`
- ▶ `IExecutionContext::setDebugListener`
- ▶ `IExecutionContext::setOutputTensorAddress`
- ▶ `IExecutionContext::setTensorDebugState`
- ▶ `IExecutionContext::updateDeviceMemorySizeForShapes`
- ▶ `IGpuAllocator::allocateAsync`
- ▶ `IGpuAllocator::deallocateAsync`
- ▶ `INetworkDefinition::addPluginV3`
- ▶ `INetworkDefinition::isDebugTensor`
- ▶ `INetworkDefinition::markDebug`
- ▶ `INetworkDefinition::unmarkDebug`

- ▶ `IPluginRegistry::acquirePluginResource`
- ▶ `IPluginRegistry::deregisterCreator`
- ▶ `IPluginRegistry::getAllCreators`
- ▶ `IPluginRegistry::getCreator`
- ▶ `IPluginRegistry::registerCreator`
- ▶ `IPluginRegistry::releasePluginResource`

## 2.4. Removed C++ APIs

Table 6. Removed C++ APIs and their Suggested Superseded API

C++ API	Superseded API
<code>BuilderFlag::kENABLE_TACTIC_HEURISTIC</code>	Builder optimization level 2
<code>BuilderFlag::kSTRICT_TYPES</code>	Use for all three flags: <ol style="list-style-type: none"> <li>1. <code>kREJECT_EMPTY_ALGORITHMS</code></li> <li>2. <code>kDIRECT_IO</code></li> <li>3. <code>kPREFER_PRECISION_CONSTRAINTS</code></li> </ol> <div style="background-color: #f0f0f0; padding: 5px; margin-top: 10px;">  Note: When removing enum members (for all enums in this list) we will be changing enumeration in the enum to have sequential numbers.           </div>
<ol style="list-style-type: none"> <li>1. <code>EngineCapability::kDEFAULT</code></li> <li>2. <code>EngineCapability::kSAFE_DLA</code></li> <li>3. <code>EngineCapability::kSAFE_GPU</code></li> </ol>	<ol style="list-style-type: none"> <li>1. <code>EngineCapability::kSTANDARD</code></li> <li>2. <code>EngineCapability::kDLA_STANDALONE</code></li> <li>3. <code>EngineCapability::kSAFETY</code></li> </ol>
<code>IAlgorithm::getAlgorithmIOInfo()</code>	<code>IAlgorithm::getAlgorithmIOInfoByIndex()</code>
<code>IAlgorithmIOInfo::getTensorFormat()</code>	The strides, data type, and vectorization information is sufficient to uniquely identify tensor formats.
<ol style="list-style-type: none"> <li>1. <code>IBuilder::buildEngineWithConfig()</code></li> <li>2. <code>IBuilder::destroy()</code></li> <li>3. <code>IBuilder::getMaxBatchSize()</code></li> <li>4. <code>IBuilder::setMaxBatchSize()</code></li> </ol>	<ol style="list-style-type: none"> <li>1. <code>IBuilder::buildSerializedNetwork()</code></li> <li>2. <code>delete ObjectName</code></li> <li>3. Implicit batch is no longer supported</li> <li>4. Implicit batch is no longer supported</li> </ol>
<ol style="list-style-type: none"> <li>1. <code>IBuilderConfig::destroy()</code></li> </ol>	<ol style="list-style-type: none"> <li>1. <code>delete ObjectName</code></li> </ol>

C++ API	Superseded API
<ul style="list-style-type: none"> <li>2. <code>IBuilderConfig::getMaxWorkspaceSize()</code></li> <li>3. <code>IBuilderConfig::getMinTimingIterations()</code></li> <li>4. <code>IBuilderConfig::setMaxWorkspaceSize()</code></li> <li>5. <code>IBuilderConfig::setMinTimingIterations()</code></li> </ul>	<ul style="list-style-type: none"> <li>2. <code>IBuilderConfig::getMemoryPoolLimit()</code> with <code>MemoryPoolType::kWORKSPACE</code></li> <li>3. <code>IBuilderConfig::getAvgTimingIterations()</code></li> <li>4. <code>IBuilderConfig::setMemoryPoolLimit()</code> with <code>MemoryPoolType::kWORKSPACE</code></li> <li>5. <code>IBuilderConfig::setAvgTimingIterations()</code></li> </ul>
<ul style="list-style-type: none"> <li>1. <code>IConvolutionLayer::getDilation()</code></li> <li>2. <code>IConvolutionLayer::getKernelSize()</code></li> <li>3. <code>IConvolutionLayer::getPadding()</code></li> <li>4. <code>IConvolutionLayer::getStride()</code></li> <li>5. <code>IConvolutionLayer::setDilation()</code></li> <li>6. <code>IConvolutionLayer::setKernelSize()</code></li> <li>7. <code>IConvolutionLayer::setPadding()</code></li> <li>8. <code>IConvolutionLayer::setStride()</code></li> </ul>	<ul style="list-style-type: none"> <li>1. <code>IConvolutionLayer::getDilationNd()</code></li> <li>2. <code>IConvolutionLayer::getKernelSizeNd()</code></li> <li>3. <code>IConvolutionLayer::getPaddingNd()</code></li> <li>4. <code>IConvolutionLayer::getStrideNd()</code></li> <li>5. <code>IConvolutionLayer::setDilationNd()</code></li> <li>6. <code>IConvolutionLayer::setKernelSizeNd()</code></li> <li>7. <code>IConvolutionLayer::setPaddingNd()</code></li> <li>8. <code>IConvolutionLayer::setStrideNd()</code></li> </ul>
<ul style="list-style-type: none"> <li>1. <code>ICudaEngine::bindingIsInput()</code></li> <li>2. <code>ICudaEngine::destroy()</code></li> <li>3. <code>ICudaEngine::getBindingBytesPerComponent()</code></li> <li>4. <code>ICudaEngine::getBindingComponentsPerElement()</code></li> <li>5. <code>ICudaEngine::getBindingDataType()</code></li> <li>6. <code>ICudaEngine::getBindingDimensions()</code></li> <li>7. <code>ICudaEngine::getBindingFormat()</code></li> <li>8. <code>ICudaEngine::getBindingFormatDesc()</code></li> <li>9. <code>ICudaEngine::getBindingIndex()</code></li> <li>10. <code>ICudaEngine::getBindingName()</code></li> <li>11. <code>ICudaEngine::getBindingVectorizedDim()</code></li> <li>12. <code>ICudaEngine::getLocation()</code></li> <li>13. <code>ICudaEngine::getMaxBatchSize()</code></li> <li>14. <code>ICudaEngine::getNbBindings()</code></li> <li>15. <code>ICudaEngine::getProfileDimensions()</code></li> <li>16. <code>ICudaEngine::getProfileShapeValues()</code></li> <li>17. <code>ICudaEngine::hasImplicitBatchDimension()</code></li> <li>18. <code>ICudaEngine::isExecutionBinding()</code></li> <li>19. <code>ICudaEngine::isShapeBinding()</code></li> </ul>	<ul style="list-style-type: none"> <li>1. <code>ICudaEngine::getTensorIOMode()</code></li> <li>2. <code>delete ObjectName</code></li> <li>3. <code>ICudaEngine::getTensorBytesPerComponent()</code></li> <li>4. <code>ICudaEngine::getTensorComponentsPerElement()</code></li> <li>5. <code>ICudaEngine::getTensorDataType()</code></li> <li>6. <code>ICudaEngine::getTensorShape()</code></li> <li>7. <code>ICudaEngine::getTensorFormat()</code></li> <li>8. <code>ICudaEngine::getTensorFormatDesc()</code></li> <li>9. <b>Name-based methods</b></li> <li>10. <b>Name-based methods</b></li> <li>11. <code>ICudaEngine::getTensorVectorizedDim()</code></li> <li>12. <code>ITensor::getLocation()</code></li> <li>13. <b>Implicit batch is no longer supported</b></li> <li>14. <code>ICudaEngine::getNbIOTensors()</code></li> <li>15. <code>ICudaEngine::getProfileShape()</code></li> <li>16. <code>ICudaEngine::getShapeValues()</code></li> <li>17. <b>Implicit batch is no longer supported</b></li> <li>18. <b>No name-based equivalent replacement</b></li> <li>19. <code>ICudaEngine::isShapeInferenceIO()</code></li> </ul>

C++ API	Superseded API
<ol style="list-style-type: none"> <li>1. IDeconvolutionLayer::getKernelSize()</li> <li>2. IDeconvolutionLayer::getPadding()</li> <li>3. IDeconvolutionLayer::getStride()</li> <li>4. IDeconvolutionLayer::setKernelSize()</li> <li>5. IDeconvolutionLayer::setPadding()</li> <li>6. IDeconvolutionLayer::setStride()</li> </ol>	<ol style="list-style-type: none"> <li>1. IDeconvolutionLayer::getKernelSizeNd()</li> <li>2. IDeconvolutionLayer::getPaddingNd()</li> <li>3. IDeconvolutionLayer::getStrideNd()</li> <li>4. IDeconvolutionLayer::setKernelSizeNd()</li> <li>5. IDeconvolutionLayer::setPaddingNd()</li> <li>6. IDeconvolutionLayer::setStrideNd()</li> </ol>
<ol style="list-style-type: none"> <li>1. IExecutionContext::destroy()</li> <li>2. IExecutionContext::enqueue()</li> <li>3. IExecutionContext::enqueueV2()</li> <li>4. IExecutionContext::execute()</li> <li>5. IExecutionContext::getBindingDimensions()</li> <li>6. IExecutionContext::getShapeBinding()</li> <li>7. IExecutionContext::getStrides()</li> <li>8. IExecutionContext::setBindingDimensions()</li> <li>9. IExecutionContext::setInputShapeBinding()</li> <li>10. IExecutionContext::setOptimizationProfile()</li> </ol>	<ol style="list-style-type: none"> <li>1. delete ObjectName</li> <li>2. IExecutionContext::enqueueV3()</li> <li>3. IExecutionContext::enqueueV3()</li> <li>4. IExecutionContext::executeV2()</li> <li>5. IExecutionContext::getTensorShape()</li> <li>6. IExecutionContext::getTensorAddress() Or getOutputTensorAddress()</li> <li>7. IExecutionContext::getTensorStrides()</li> <li>8. IExecutionContext::setInputShape()</li> <li>9. IExecutionContext::setInputTensorAddress() Or setTensorAddress()</li> <li>10. IExecutionContext::setOptimizationProfileAsync()</li> </ol>
IFullyConnectedLayer	IMatrixMultiplyLayer
IGpuAllocator::free()	IGpuAllocator::deallocate()
IHostMemory::destroy()	delete ObjectName
<ol style="list-style-type: none"> <li>1. INetworkDefinition::addConvolution()</li> <li>2. INetworkDefinition::addDeconvolution()</li> <li>3. INetworkDefinition::addFullyConnected()</li> <li>4. INetworkDefinition::addPadding()</li> <li>5. INetworkDefinition::addPooling()</li> <li>6. INetworkDefinition::addRNNv2()</li> <li>7. INetworkDefinition::destroy()</li> <li>8. INetworkDefinition::hasExplicitPrecision()</li> <li>9. INetworkDefinition::hasImplicitBatchDimension()</li> </ol>	<ol style="list-style-type: none"> <li>1. INetworkDefinition::addConvolutionNd()</li> <li>2. INetworkDefinition::addDeconvolutionNd()</li> <li>3. INetworkDefinition::addMatrixMultiply()</li> <li>4. INetworkDefinition::addPaddingNd()</li> <li>5. INetworkDefinition::addPoolingNd()</li> <li>6. INetworkDefinition::addLoop()</li> <li>7. delete ObjectName</li> <li>8. Explicit precision support is removed in 10.0</li> <li>9. Implicit batch support is removed</li> </ol>
IOnnxConfig::destroy()	delete ObjectName
<ol style="list-style-type: none"> <li>1. IPaddingLayer::getPostPadding()</li> </ol>	<ol style="list-style-type: none"> <li>1. IPaddingLayer::getPostPaddingNd()</li> </ol>

C++ API	Superseded API
<ol style="list-style-type: none"> <li>IPaddingLayer::getPrePadding()</li> <li>IPaddingLayer::setPostPadding()</li> <li>IPaddingLayer::setPrePadding()</li> </ol>	<ol style="list-style-type: none"> <li>IPaddingLayer::getPrePaddingNd()</li> <li>IPaddingLayer::setPostPaddingNd()</li> <li>IPaddingLayer::setPrePaddingNd()</li> </ol>
<ol style="list-style-type: none"> <li>IPoolingLayer::getPadding()</li> <li>IPoolingLayer::getStride()</li> <li>IPoolingLayer::getWindowSize()</li> <li>IPoolingLayer::setPadding()</li> <li>IPoolingLayer::setStride()</li> <li>IPoolingLayer::setWindowSize()</li> </ol>	<ol style="list-style-type: none"> <li>IPoolingLayer::getPaddingNd()</li> <li>IPoolingLayer::getStrideNd()</li> <li>IPoolingLayer::getWindowSizeNd()</li> <li>IPoolingLayer::setPaddingNd()</li> <li>IPoolingLayer::setStrideNd()</li> <li>IPoolingLayer::setWindowSizeNd()</li> </ol>
IRefitter::destroy()	delete ObjectName
<ol style="list-style-type: none"> <li>IResizeLayer::getAlignCorners()</li> <li>IResizeLayer::setAlignCorners()</li> </ol>	<ol style="list-style-type: none"> <li>IResizeLayer::getAlignCornersNd()</li> <li>IResizeLayer::setAlignCornersNd()</li> </ol>
<ol style="list-style-type: none"> <li>IRuntime::deserializeCudaEngine(void const* blob, std::size_t size, IPluginFactory* pluginFactory)</li> <li>IRuntime::destroy()</li> </ol>	<ol style="list-style-type: none"> <li>Use deserializeCudaEngine with two parameters</li> <li>delete ObjectName</li> </ol>
IRNNv2Layer	ILoop
kNV_TENSORRT_VERSION_IMPL	<pre>#define NV_TENSORRT_VERSION_INT(major, minor, patch) ((major) *10000L + (minor) *100L + (patch) *1L)</pre> <div style="background-color: #f0f0f0; padding: 5px; margin-top: 10px;">  Note: TensorRT version encoding was changed to accommodate a two digit minor version. </div>
<ol style="list-style-type: none"> <li>NetworkDefinitionCreationFlag::kEXPLICIT_BATCH</li> <li>NetworkDefinitionCreationFlag::kEXPLICIT_PRECISION</li> </ol>	Support is removed in 10.0
<ol style="list-style-type: none"> <li>NV_TENSORRT_SONAME_MAJOR</li> <li>NV_TENSORRT_SONAME_MINOR</li> <li>NV_TENSORRT_SONAME_PATCH</li> </ol>	<ol style="list-style-type: none"> <li>NV_TENSORRT_MAJOR</li> <li>NV_TENSORRT_MINOR</li> <li>NV_TENSORRT_PATCH</li> </ol>
<ol style="list-style-type: none"> <li>PaddingMode::kCAFFE_ROUND_DOWN</li> <li>PaddingMode::kCAFFE_ROUND_UP</li> </ol>	Caffe is not supported since 9.0

C++ API	Superseded API
<ol style="list-style-type: none"> <li>1. <code>PreviewFeature::kDISABLE_EXTERNAL_TACTICS</code></li> <li>2. <code>PreviewFeature::kFASTER_DYNAMIC_SHAPES</code></li> </ol>	<ol style="list-style-type: none"> <li>1. External tactics are always disabled for core code</li> <li>2. This flag is on by default</li> </ol>
<ol style="list-style-type: none"> <li>1. <code>ProfilingVerbosity::kDEFAULT</code></li> <li>2. <code>ProfilingVerbosity::kVERBOSE</code></li> </ol>	<ol style="list-style-type: none"> <li>1. <code>ProfilingVerbosity::kLAYER_NAMES_ONLY</code></li> <li>2. <code>ProfilingVerbosity::kDETAILED</code></li> </ol>
<code>ResizeMode</code>	Use <code>InterpolationMode</code> , alias is removed
<ol style="list-style-type: none"> <li>1. <code>RNNDirection</code></li> <li>2. <code>RNNGateType</code></li> <li>3. <code>RNNInputMode</code></li> <li>4. <code>RNNOperation</code></li> </ol>	RNN related data structures are removed
<code>SampleMode::kDEFAULT</code>	<code>SampleMode::kSTRICT_BOUNDS</code>
<code>SliceMode</code>	Use <code>SampleMode</code> , alias is removed

## 2.5. Removed C++ Plugins

Table 7. Removed C++ Plugins and their Suggested Superseded Plugin

C++ Plugin	Superseded Plugin
1. <code>createAnchorGeneratorPlugin()</code>	1. <code>GridAnchorPluginCreator::createPlugin()</code>
2. <code>createBatchedNMSPlugin()</code>	2. <code>BatchedNMSPluginCreator::createPlugin()</code>
3. <code>createInstanceNormalizationPlugin()</code>	3. <code>InstanceNormalizationPluginCreator::createPlugin()</code>
4. <code>createNMSPlugin()</code>	4. <code>NMSPluginCreator::createPlugin()</code>
5. <code>createNormalizePlugin()</code>	5. <code>NormalizePluginCreator::createPlugin()</code>
6. <code>createPriorBoxPlugin()</code>	6. <code>PriorBoxPluginCreator::createPlugin()</code>
7. <code>createRegionPlugin()</code>	7. <code>RegionPluginCreator::createPlugin()</code>
8. <code>createReorgPlugin()</code>	8. <code>ReorgPluginCreator::createPlugin()</code>
9. <code>createRPNROIPlugin()</code>	9. <code>RPROIPluginCreator::createPlugin()</code>
10. <code>createSplitPlugin()</code>	10. <code>INetworkDefinition::addSlice()</code>
<code>struct Quadruple</code>	Related plugins are removed



## 2.6. Removed Safety C++ APIs

Table 8. Removed Safety C++ APIs and their Suggested Superseded Safety API

Safety C++ API	Superseded Safety API
1. <code>safe::ICudaEngine::bindingIsInput()</code>	1. <code>safe::ICudaEngine::tensorIOMode()</code>
2. <code>safe::ICudaEngine::getBindingBytesPerComponent()</code>	2. <code>safe::ICudaEngine::getTensorBytesPerComponent()</code>
3. <code>safe::ICudaEngine::getBindingComponentsPerElement()</code>	3. <code>safe::ICudaEngine::getTensorComponentsPerElement()</code>
4. <code>safe::ICudaEngine::getBindingDataType()</code>	4. <code>safe::ICudaEngine::getTensorDataType()</code>
5. <code>safe::ICudaEngine::getBindingDimensions()</code>	5. <code>safe::ICudaEngine::getTensorShape()</code>
6. <code>safe::ICudaEngine::getBindingIndex()</code>	6. <code>safe::name-based methods</code>
7. <code>safe::ICudaEngine::getBindingName()</code>	7. <code>safe::name-based methods</code>
8. <code>safe::ICudaEngine::getBindingVectorizedDim()</code>	8. <code>safe::ICudaEngine::getTensorVectorizedDim()</code>
9. <code>safe::ICudaEngine::getNbBindings()</code>	9. <code>safe::ICudaEngine::getNbIOTensors()</code>
10. <code>safe::ICudaEngine::getTensorFormat()</code>	10. <code>safe::ICudaEngine::getBindingFormat()</code>
1. <code>safe::IExecutionContext::enqueueV2()</code>	1. <code>safe::IExecutionContext::enqueueV3()</code>
2. <code>safe::IExecutionContext::getStrides()</code>	2. <code>safe::IExecutionContext::getTensorStrides()</code>

# Chapter 3. trtexec

## 3.1. trtexec Flag Changes

Table 9. Changes to flag workspace and minTiming

TensorRT 8.x	TensorRT 10.0
<pre>trtexec \   --onnx=/path/to/model.onnx \   --saveEngine=/path/to/engine.trt \   --optShapes=input:\$INPUT_SHAPE \   --avgTiming=1 \   --workspace=1024 \   --minTiming=1</pre>	<pre>trtexec \   --onnx=/path/to/model.onnx \   --saveEngine=/path/to/engine.trt \   --optShapes=input:\$INPUT_SHAPE \   --avgTiming=1 \   --memPoolSize=workspace:1024</pre>

## 3.2. Removed trtexec Flags

Table 10. Removed trtexec Flags and their Suggested Superseded Flag

trtexec Flag	Superseded Flag
<code>--minTiming</code>	<code>avgTiming</code>
<code>--preview=features options:</code> <ul style="list-style-type: none"><li>▶ <code>disableExternalTacticSourcesForCore0805</code></li><li>▶ <code>fasterDynamicShapes0805</code></li></ul>	N/A
<code>--workspace=N</code>	<code>--memPoolSize=poolspec</code>

## 3.3. Deprecated trtexec Flags

Table 11. Deprecated trtexec Flags

<b>trtexec Flag</b>
<code>--buildOnly</code>
<code>--explicitPrecision</code>
<code>--heuristic</code>
<code>--nvtxMode</code>

## 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, Hopper, JetPack, Jetson AGX Xavier, Jetson Nano, 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

© 2024-2024 NVIDIA Corporation & affiliates. All rights reserved.

