



# CUDNN LIBRARY

DU-06702-001\_v6.5 | August 2014

## User Guide





# Chapter 1.

## INTRODUCTION

NVIDIA® cuDNN is a GPU-accelerated library of primitives for deep neural networks. It provides highly tuned implementations of routines arising frequently in DNN applications:

- ▶ Convolution forward and backward, including cross-correlation
- ▶ Pooling forward and backward
- ▶ Softmax forward and backward
- ▶ Neuron activations forward and backward:
  - ▶ Rectified linear (ReLU)
  - ▶ Sigmoid
  - ▶ Hyperbolic tangent (TANH)
- ▶ Tensor transformation functions

cuDNN's convolution routines aim for performance competitive with the fastest GEMM (matrix multiply) based implementations of such routines while using significantly less memory.

cuDNN features customizable data layouts, supporting flexible dimension ordering, striding, and subregions for the 4D tensors used as inputs and outputs to all of its routines. This flexibility allows easy integration into any neural network implementation and avoids the input/output transposition steps sometimes necessary with GEMM-based convolutions.

cuDNN offers a context-based API that allows for easy multithreading and (optional) interoperability with CUDA streams.

# Chapter 2.

## GENERAL DESCRIPTION

### 2.1. Programming Model

The cuDNN Library exposes a Host API but assumes that for operations using the GPU, the necessary data is directly accessible from the device.

An application using cuDNN must initialize a handle to the library context by calling **cudaDnnCreate()**. This handle is explicitly passed to every subsequent library function that operates on GPU data. Once the application finishes using cuDNN, it can release the resources associated with the library handle using **cudaDnnDestroy()**. This approach allows the user to explicitly control the library's functioning when using multiple host threads, GPUs and CUDA Streams. For example, an application can use **cudaSetDevice()** to associate different devices with different host threads and in each of those host threads, use a unique cuDNN handle which directs library calls to the device associated with it. cuDNN library calls made with different handles will thus, automatically run on different devices. The device associated with a particular cuDNN context is assumed to remain unchanged between the corresponding **cudaDnnCreate()** and **cudaDnnDestroy()** calls. In order for the cuDNN library to use a different device within the same host thread, the application must set the new device to be used by calling **cudaSetDevice()** and then create another cuDNN context, which will be associated with the new device, by calling **cudaDnnCreate()**.

### 2.2. Thread Safety

The library is thread safe and its functions can be called from multiple host threads, even with the same handle. When sharing a handle across host threads, extreme care needs to be taken to ensure that any changes to the handle configuration in one thread do not adversely affect cuDNN function calls in others. This is especially true for the destruction of the handle. It is not recommended that multiple threads share the same cuDNN handle.

## 2.3. Reproducibility

By design, most of cuDNN's routines from a given version generate the same bit-wise results at every run when executed on GPUs with the same architecture and the same number of SMs. However, bit-wise reproducibility is not guaranteed across versions, as the implementation of a given routine may change. With the current release, the following routines do not guarantee reproducibility because they use atomic add operations:

- ▶ `cudaDnnConvolutionBackwardFilter`
- ▶ `cudaDnnConvolutionBackwardData`

## 2.4. Scaling parameters `alpha` and `beta`

Many cuDNN routines like `cudaDnnConvolutionForward` support the use of scalar parameters `alpha` and `beta` to scale the input and the output tensors respectively. These parameters, passed by reference, should be of the same type as the tensor they are applied to and should reside on the Host. When `beta` is equal to zero, the output tensor can be provided as uninitialized : it will not be read and any potential `Nan` will have no effect.

## 2.5. Requirements

cuDNN supports NVIDIA GPUs of compute capability 3.0 and higher and requires an NVIDIA Driver compatible with CUDA Toolkit 6.5.

# Chapter 3.

## CUDNN DATATYPES REFERENCE

This chapter describes all the types and enums of the cuDNN library API.

### 3.1. cudnnHandle\_t

**cudnnHandle\_t** is a pointer to an opaque structure holding the cuDNN library context. The cuDNN library context must be created using **cudnnCreate()** and the returned handle must be passed to all subsequent library function calls. The context should be destroyed at the end using **cudnnDestroy()**. The context is associated with only one GPU device, the current device at the time of the call to **cudnnCreate()**. However multiple contexts can be created on the same GPU device.

### 3.2. cudnnStatus\_t

**cudnnStatus\_t** is an enumerated type used for function status returns. All cuDNN library functions return their status, which can be one of the following values:

Value	Meaning
CUDNN_STATUS_SUCCESS	The operation completed successfully.
CUDNN_STATUS_NOT_INITIALIZED	The cuDNN library was not initialized properly. This error is usually returned when a call to <b>cudnnCreate()</b> fails or when <b>cudnnCreate()</b> has not been called prior to calling another cuDNN routine. In the former case, it is usually due to an error in the CUDA Runtime API called by <b>cudnnCreate()</b> or by an error in the hardware setup.
CUDNN_STATUS_ALLOC_FAILED	Resource allocation failed inside the cuDNN library. This is usually caused by an internal <b>cudaMalloc()</b> failure.  To correct: prior to the function call, deallocate previously allocated memory as much as possible.

Value	Meaning
<code>CUDNN_STATUS_BAD_PARAM</code>	<p>An incorrect value or parameter was passed to the function.</p> <p>To correct: ensure that all the parameters being passed have valid values.</p>
<code>CUDNN_STATUS_ARCH_MISMATCH</code>	<p>The function requires a feature absent from the current GPU device. Note that cuDNN only supports devices with compute capabilities greater than or equal to 3.0.</p> <p>To correct: compile and run the application on a device with appropriate compute capability.</p>
<code>CUDNN_STATUS_MAPPING_ERROR</code>	<p>An access to GPU memory space failed, which is usually caused by a failure to bind a texture.</p> <p>To correct: prior to the function call, unbind any previously bound textures.</p> <p>Otherwise, this may indicate an internal error/bug in the library.</p>
<code>CUDNN_STATUS_EXECUTION_FAILED</code>	<p>The GPU program failed to execute. This is usually caused by a failure to launch some cuDNN kernel on the GPU, which can occur for multiple reasons.</p> <p>To correct: check that the hardware, an appropriate version of the driver, and the cuDNN library are correctly installed.</p> <p>Otherwise, this may indicate a internal error/bug in the library.</p>
<code>CUDNN_STATUS_INTERNAL_ERROR</code>	An internal cuDNN operation failed.
<code>CUDNN_STATUS_NOT_SUPPORTED</code>	The functionality requested is not presently supported by cuDNN.
<code>CUDNN_STATUS_LICENSE_ERROR</code>	<p>The functionality requested requires some license and an error was detected when trying to check the current licensing. This error can happen if the license is not present or is expired or if the environment variable <code>NVIDIA_LICENSE_FILE</code> is not set properly.</p>

### 3.3. `cudaTensorDescriptor_t`

`cudaCreateTensorDescriptor_t` is a pointer to an opaque structure holding the description of a generic n-D dataset. `cudaCreateTensorDescriptor()` is used to create one instance, and one of the routines `cudaSetTensorNdDescriptor()`, `cudaSetTensor4dDescriptor()` or `cudaSetTensor4dDescriptorEx()` must be used to initialize this instance.

### 3.4. cudnnFilterDescriptor\_t

**cudnnFilterDescriptor\_t** is a pointer to an opaque structure holding the description of a filter dataset. **cudnnCreateFilterDescriptor()** is used to create one instance, and **cudnnSetFilterDescriptor()** must be used to initialize this instance.

### 3.5. cudnnConvolutionDescriptor\_t

**cudnnConvolutionDescriptor\_t** is a pointer to an opaque structure holding the description of a convolution operation. **cudnnCreateConvolutionDescriptor()** is used to create one instance, and **cudnnSetConvolutionNdDescriptor()** or **cudnnSetConvolution2dDescriptor()** must be used to initialize this instance.

### 3.6. cudnnPoolingDescriptor\_t

**cudnnPoolingDescriptor\_t** is a pointer to an opaque structure holding the description of a pooling operation. **cudnnCreatePoolingDescriptor()** is used to create one instance, and **cudnnSetPoolingNdDescriptor()** or **cudnnSetPooling2dDescriptor()** must be used to initialize this instance.

### 3.7. cudnnDataType\_t

**cudnnDataType\_t** is an enumerated type indicating the data type to which a tensor descriptor or filter descriptor refers.

Value	Meaning
CUDNN_DATA_FLOAT	The data is 32-bit single-precision floating point (float).
CUDNN_DATA_DOUBLE	The data is 64-bit double-precision floating point (double).

### 3.8. cudnnTensorFormat\_t

**cudnnTensorFormat\_t** is an enumerated type used by **cudnnSetTensor4dDescriptor()** to create a tensor with a pre-defined layout.

Value	Meaning
CUDNN_TENSOR_NCHW	This tensor format specifies that the data is laid out in the following order: image, features map, rows, columns. The strides are implicitly defined in such a way that the data are contiguous in memory with no padding between images, feature maps, rows, and columns; the columns are the



Value	Meaning
	inner dimension and the images are the outermost dimension.
<code>CUDNN_TENSOR_NHWC</code>	This tensor format specifies that the data is laid out in the following order: image, rows, columns, features maps. The strides are implicitly defined in such a way that the data are contiguous in memory with no padding between images, rows, columns, and features maps; the feature maps are the inner dimension and the images are the outermost dimension.

### 3.9. `cudaAddMode_t`

`cudaAddMode_t` is an enumerated type used by `cudaAddTensor()` to specify how a bias tensor is added to an input/output tensor.

Value	Meaning
<code>CUDNN_ADD_IMAGE</code> or <code>CUDNN_ADD_SAME_HW</code>	In this mode, the bias tensor is defined as one image with one feature map. This image will be added to every feature map of every image of the input/output tensor.
<code>CUDNN_ADD_FEATURE_MAP</code> or <code>CUDNN_ADD_SAME_CHW</code>	In this mode, the bias tensor is defined as one image with multiple feature maps. This image will be added to every image of the input/output tensor.
<code>CUDNN_ADD_SAME_C</code>	In this mode, the bias tensor is defined as one image with multiple feature maps of dimension 1x1; it can be seen as an vector of feature maps. Each feature map of the bias tensor will be added to the corresponding feature map of all height-by-width pixels of every image of the input/output tensor.
<code>CUDNN_ADD_FULL_TENSOR</code>	In this mode, the bias tensor has the same dimensions as the input/output tensor. It will be added point-wise to the input/output tensor.

### 3.10. `cudaConvolutionMode_t`

`cudaConvolutionMode_t` is an enumerated type used by `cudaSetConvolutionDescriptor()` to configure a convolution descriptor. The filter used for the convolution can be applied in two different ways, corresponding mathematically to a convolution or to a cross-correlation. (A cross-correlation is equivalent to a convolution with its filter rotated by 180 degrees.)

Value	Meaning
CUDNN_CONVOLUTION	In this mode, a convolution operation will be done when applying the filter to the images.
CUDNN_CROSS_CORRELATION	In this mode, a cross-correlation operation will be done when applying the filter to the images.

### 3.11. cudnnConvolutionFwdPreference\_t

**cudnnConvolutionFwdPreference\_t** is an enumerated type used by **cudnnGetConvolutionForwardAlgorithm()** to help the choice of the algorithm used for the forward convolution.

Value	Meaning
CUDNN_CONVOLUTION_FWD_NO_WORKSPACE	In this configuration, the routine <b>cudnnGetConvolutionForwardAlgorithm()</b> is guaranteed to return an algorithm that does not require any extra workspace to be provided by the user.
CUDNN_CONVOLUTION_FWD_PREFER_FASTEST	In this configuration, the routine <b>cudnnGetConvolutionForwardAlgorithm()</b> will return the fastest algorithm regardless how much workspace is needed to execute it.
CUDNN_CONVOLUTION_FWD_SPECIFY_WORKSPACE_LIMIT	In this configuration, the routine <b>cudnnGetConvolutionForwardAlgorithm()</b> will return the fastest algorithm that fits within the memory limit that the user provided.

### 3.12. cudnnConvolutionFwdAlgo\_t

**cudnnConvolutionFwdAlgo\_t** is an enumerated type that exposes the different algorithm available to execute the forward convolution operation.

Value	Meaning
CUDNN_CONVOLUTION_FWD_ALGO_IMPLICIT_GEMM	This algorithm expresses the convolution as a matrix product without actually explicitly form the matrix that holds the input tensor data.
CUDNN_CONVOLUTION_FWD_ALGO_IMPLICIT_PRECOMPUTED_GEMM	This algorithm expresses the convolution as a matrix product without actually explicitly form the matrix that holds the input tensor data, but still needs some memory workspace to precompute some indices in order to facilitate the implicit construction of the matrix that holds the input tensor data
CUDNN_CONVOLUTION_FWD_ALGO_GEMM	This algorithm expresses the convolution as an explicit matrix product. A significant memory

Value	Meaning
	workspace is needed to store the matrix that holds the input tensor data.
<code>CUDNN_CONVOLUTION_FWD_ALGO_DIRECT</code>	This algorithm expresses the convolution as a direct convolution (e.g without implicitly or explicitly doing a matrix multiplication).

### 3.13. `cudannSoftmaxAlgorithm_t`

`cudannSoftmaxAlgorithm_t` is used to select an implementation of the softmax function used in `cudannSoftmaxForward()` and `cudannSoftmaxBackward()`.

Value	Meaning
<code>CUDNN_SOFTMAX_FAST</code>	This implementation applies the straightforward softmax operation.
<code>CUDNN_SOFTMAX_ACCURATE</code>	This implementation scales each point of the softmax input domain by its maximum value to avoid potential floating point overflows in the softmax evaluation.

### 3.14. `cudannSoftmaxMode_t`

`cudannSoftmaxMode_t` is used to select over which data the `cudannSoftmaxForward()` and `cudannSoftmaxBackward()` are computing their results.

Value	Meaning
<code>CUDNN_SOFTMAX_MODE_INSTANCE</code>	The softmax operation is computed per image (N) across the dimensions C,H,W.
<code>CUDNN_SOFTMAX_MODE_CHANNEL</code>	The softmax operation is computed per spatial location (H,W) per image (N) across the dimension C.

### 3.15. `cudannPoolingMode_t`

`cudannPoolingMode_t` is an enumerated type passed to `cudannSetPoolingDescriptor()` to select the pooling method to be used by `cudannPoolingForward()` and `cudannPoolingBackward()`.

Value	Meaning
<code>CUDNN_POOLING_MAX</code>	The maximum value inside the pooling window will be used.
<code>CUDNN_POOLING_AVERAGE_COUNT_INCLUDE_PADDING</code>	The values inside the pooling window will be averaged. The number of padded values will be

Value	Meaning
	taken into account when computing the average value
CUDNN_POOLING_AVERAGE_COUNT_EXCLUDE_PADDING	The values inside the pooling window will be averaged. The number of padded values will not be taken into account when computing the average value

## 3.16. cudnnActivationMode\_t

**cudnnActivationMode\_t** is an enumerated type used to select the neuron activation function used in **cudnnActivationForward()** and **cudnnActivationBackward()**.

Value	Meaning
CUDNN_ACTIVATION_SIGMOID	Selects the sigmoid function.
CUDNN_ACTIVATION_RELU	Selects the rectified linear function.
CUDNN_ACTIVATION_TANH	Selects the hyperbolic tangent function.

## 3.17. cudnnDataType\_t

**cudnnDataType\_t** is an enumerated type indicating the data type to which a tensor descriptor or filter descriptor refers.

Value	Meaning
CUDNN_DATA_FLOAT	The data is 32-bit single-precision floating point ( <b>float</b> ).
CUDNN_DATA_DOUBLE	The data is 64-bit double-precision floating point ( <b>double</b> ).

# Chapter 4.

## CUDNN API REFERENCE

This chapter describes the API of all the routines of the cuDNN library.

### 4.1. cudnnGetVersion

```
size_t cudnnGetVersion()
```

This function returns the version number of the cuDNN Library. It returns the **CUDNN\_VERSION** define present in the cudnn.h header file. Starting with release R2, the routine can be used to identify dynamically the current cuDNN Library used by the application. The define **CUDNN\_VERSION** can be used to have the same application linked against different cuDNN versions using conditional compilation statements.

### 4.2. cudnnGetErrorString

```
const char * cudnnGetErrorString(cudnnStatus_t status)
```

This function returns a human-readable character string describing the **cudnnStatus\_t** enumerate passed as input parameter.

### 4.3. cudnnCreate

```
cudnnStatus_t cudnnCreate(cudnnHandle_t *handle)
```

This function initializes the cuDNN library and creates a handle to an opaque structure holding the cuDNN library context. It allocates hardware resources on the host and device and must be called prior to making any other cuDNN library calls. The cuDNN library context is tied to the current CUDA device. To use the library on multiple devices, one cuDNN handle needs to be created for each device. For a given device, multiple cuDNN handles with different configurations (e.g., different current CUDA streams) may be created. Because **cudnnCreate** allocates some internal resources, the release of those resources by calling **cudnnDestroy** will implicitly call **cudaDeviceSynchronize**; therefore, the recommended best practice is to call **cudnnCreate/cudnnDestroy** outside of performance-critical code paths. For multithreaded applications that use the same device from different threads, the

recommended programming model is to create one (or a few, as is convenient) cuDNN handle(s) per thread and use that cuDNN handle for the entire life of the thread.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The initialization succeeded.
<code>CUDNN_STATUS_NOT_INITIALIZED</code>	CUDA Runtime API initialization failed.
<code>CUDNN_STATUS_ALLOC_FAILED</code>	The resources could not be allocated.

## 4.4. cudnnDestroy

```
cudaStatus_t cudnnDestroy(cudaHandle_t handle)
```

This function releases hardware resources used by the cuDNN library. This function is usually the last call with a particular handle to the cuDNN library. Because **cudnnCreate** allocates some internal resources, the release of those resources by calling **cudnnDestroy** will implicitly call **cudaDeviceSynchronize**; therefore, the recommended best practice is to call **cudnnCreate/cudnnDestroy** outside of performance-critical code paths.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The cuDNN context destruction was successful.
<code>CUDNN_STATUS_NOT_INITIALIZED</code>	The library was not initialized.

## 4.5. cudnnSetStream

```
cudaStatus_t cudnnSetStream(cudaHandle_t handle, cudaStream_t streamId)
```

This function sets the cuDNN library stream, which will be used to execute all subsequent calls to the cuDNN library functions with that particular handle. If the cuDNN library stream is not set, all kernels use the default (**NULL**) stream. In particular, this routine can be used to change the stream between kernel launches and then to reset the cuDNN library stream back to **NULL**.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The stream was set successfully.

## 4.6. cudnnGetStream

```
cudaStatus_t cudnnGetStream(cudaHandle_t handle, cudaStream_t *streamId)
```

This function gets the cuDNN library stream, which is being used to execute all calls to the cuDNN library functions. If the cuDNN library stream is not set, all kernels use the *default* **NULL** stream.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The stream was returned successfully.

## 4.7. cudnnCreateTensorDescriptor

```

cudnnStatus_t cudnnCreateTensorDescriptor(cudnnTensorDescriptor_t *tensorDesc)

```

This function creates a generic Tensor descriptor object by allocating the memory needed to hold its opaque structure.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was created successfully.
CUDNN_STATUS_ALLOC_FAILED	The resources could not be allocated.

## 4.8. cudnnSetTensor4dDescriptor

```

cudnnStatus_t
cudnnSetTensor4dDescriptor( cudnnTensorDescriptor_t tensorDesc,
                           cudnnTensorFormat_t format,
                           cudnnDataType_t dataType,
                           int n,
                           int c,
                           int h,
                           int w )

```

This function initializes a previously created generic Tensor descriptor object into a 4D tensor. The strides of the four dimensions are inferred from the format parameter and set in such a way that the data is contiguous in memory with no padding between dimensions.



The total size of a tensor including the potential padding between dimensions is limited to 2 Giga-elements of type `datatype`.

Param	In/out	Meaning
tensorDesc	input/output	Handle to a previously created tensor descriptor.
format	input	Type of format.
datatype	input	Data type.
n	input	Number of images.
c	input	Number of feature maps per image.
h	input	Height of each feature map.
w	input	Width of each feature map.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was set successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the parameters <code>n</code> , <code>c</code> , <code>h</code> , <code>w</code> was negative or <code>format</code> has an invalid enumerant value or <code>dataType</code> has an invalid enumerant value.
CUDNN_STATUS_NOT_SUPPORTED	The total size of the tensor descriptor exceeds the maximim limit of 2 Giga-elements.

## 4.9. cudnnSetTensor4dDescriptorEx

```

cudnnStatus_t
cudnnSetTensor4dDescriptorEx( cudnnTensorDescriptor_t tensorDesc,
                             cudnnDataType_t dataType,
                             int n,
                             int c,
                             int h,
                             int w,
                             int nStride,
                             int cStride,
                             int hStride,
                             int wStride )

```

This function initializes a previously created generic Tensor descriptor object into a 4D tensor, similarly to **cudnnSetTensor4dDescriptor** but with the strides explicitly passed as parameters. This can be used to lay out the 4D tensor in any order or simply to define gaps between dimensions.



At present, some cuDNN routines have limited support for strides; Those routines will return CUDNN\_STATUS\_NOT\_SUPPORTED if a Tensor4D object with an unsupported stride is used. **cudnnTransformTensor** can be used to convert the data to a supported layout.



The total size of a tensor including the potential padding between dimensions is limited to 2 Giga-elements of type `dataType`.

Param	In/out	Meaning
tensorDesc	input/ output	Handle to a previously created tensor descriptor.
datatype	input	Data type.
n	input	Number of images.
c	input	Number of feature maps per image.
h	input	Height of each feature map.
w	input	Width of each feature map.
nStride	input	Stride between two consecutive images.



Param	In/out	Meaning
cStride	input	Stride between two consecutive feature maps.
hStride	input	Stride between two consecutive rows.
wStride	input	Stride between two consecutive columns.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was set successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the parameters <code>n</code> , <code>c</code> , <code>h</code> , <code>w</code> or <code>nStride</code> , <code>cStride</code> , <code>hStride</code> , <code>wStride</code> is negative or <code>dataType</code> has an invalid enumerant value.
CUDNN_STATUS_NOT_SUPPORTED	The total size of the tensor descriptor exceeds the maximim limit of 2 Giga-elements.

## 4.10. cudnnGetTensor4dDescriptor

```

cudnnStatus_t
cudnnGetTensor4dDescriptor( cudnnTensorDescriptor_t tensorDesc,
                           cudnnDataType_t *dataType,
                           int *n,
                           int *c,
                           int *h,
                           int *w,
                           int *nStride,
                           int *cStride,
                           int *hStride,
                           int *wStride )

```

This function queries the parameters of the previously initialized Tensor4D descriptor object.

Param	In/out	Meaning
tensorDesc	input	Handle to a previously insitialized tensor descriptor.
datatype	output	Data type.
n	output	Number of images.
c	output	Number of feature maps per image.
h	output	Height of each feature map.
w	output	Width of each feature map.
nStride	output	Stride between two consecutive images.
cStride	output	Stride between two consecutive feature maps.
hStride	output	Stride between two consecutive rows.
wStride	output	Stride between two consecutive columns.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The operation succeeded.

## 4.11. cudnnSetTensorNdDescriptor

```

cudnnStatus_t
cudnnSetTensorNdDescriptor( cudnnTensorDescriptor_t  tensorDesc,

                           cudnnDataType_t dataType,
                           int nbDims,
                           int dimA[],
                           int strideA[])

```

This function initializes a previously created generic Tensor descriptor object.



The total size of a tensor including the potential padding between dimensions is limited to 2 Giga-elements of type `dataType`.

Param	In/out	Meaning
tensorDesc	input/ output	Handle to a previously created tensor descriptor.
dataType	input	Data type.
nbDims	input	Dimension of the tensor.
dimA	input	Array of dimension <code>nbDims</code> that contain the size of the tensor for every dimension.
strideA	input	Array of dimension <code>nbDims</code> that contain the stride of the tensor for every dimension.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The object was set successfully.
<code>CUDNN_STATUS_BAD_PARAM</code>	At least one of the elements of the array <code>dimA</code> was negative or <code>dataType</code> has an invalid enumerant value.
<code>CUDNN_STATUS_NOT_SUPPORTED</code>	the parameter <code>nbDims</code> exceeds the maximum supported dimension or the total size of the tensor descriptor exceeds the maximim limit of 2 Giga-elements.

## 4.12. cudnnGetTensorNdDescriptor

```

cudnnStatus_t
cudnnGetTensorNdDescriptor( const cudnnTensorDescriptor_t  tensorDesc,

                           int nbDimsRequested,
                           cudnnDataType_t *dataType,
                           int *nbDims,
                           int dimA[],
                           int strideA[])

```

This function requires a previously initialized generic Tensor descriptor object.

Param	In/out	Meaning
tensorDesc	input	Handle to a previously initialized tensor descriptor.
nbDimsRequested	input	Dimension of the expected tensor descriptor. It is also the minimum size of the arrays <b>dimA</b> and <b>strideA</b> in order to be able to hold the results
datatype	output	Data type.
nbDims	output	Actual dimension of the tensor.
dimA	output	Array of dimension of at least <b>nbDimsRequested</b> that will be filled with the size parameters from the provided tensor descriptor.
strideA	input	Array of dimension of at least <b>nbDimsRequested</b> that will be filled with the stride parameters from the provided tensor descriptor.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was set successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the elements of the array <b>dimA</b> was negative or <b>dataType</b> has an invalid enumerant value.
CUDNN_STATUS_NOT_SUPPORTED	the parameter <b>nbDims</b> exceeds the maximum supported dimension

## 4.13. cudnnDestroyTensorDescriptor

```

cudnnStatus_t cudnnDestroyTensorDescriptor(cudnnTensorDescriptor_t tensorDesc)

```

This function destroys a previously created Tensor descriptor object.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was destroyed successfully.

## 4.14. cudnnTransformTensor

```

cudnnStatus_t
cudnnTransformTensor( cudnnHandle_t      handle,
                     const void          *alpha,
                     const cudnnTensorDescriptor_t srcDesc,
                     const void          *srcData,
                     const void          *beta,
                     const cudnnTensorDescriptor_t destDesc,
                     void                *destData )

```

This function copies the scaled data from one tensor to another tensor with a different layout. Those descriptors need to have the same dimensions but not necessarily the same strides. The input and output tensors must not overlap in any way (i.e., tensors cannot be transformed in place). This function can be used to convert a tensor with an unsupported format to a supported one.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
alpha	input	Scalar factor to be applied to every element of the input tensor before it is added to the output tensor.
srcDesc	input	Handle to a previously initialized tensor descriptor.
srcData	input	Pointer to data of the tensor described by the <code>srcDesc</code> descriptor.
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the operation. Note that if <code>beta</code> is zero, the output is not read and can contain any uninitialized data (including Nan numbers).
destDesc	input	Handle to a previously initialized tensor descriptor.
destData	output	Pointer to data of the tensor described by the <code>destDesc</code> descriptor.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The function launched successfully.
<code>CUDNN_STATUS_BAD_PARAM</code>	The dimensions <code>n, c, h, w</code> or the <code>dataType</code> of the two tensor descriptors are different.
<code>CUDNN_STATUS_EXECUTION_FAILED</code>	The function failed to launch on the GPU.

## 4.15. cudnnAddTensor

```

cudnnStatus_t
cudnnAddTensor(  cudnnHandle_t      handle,
                  cudnnAddMode_t    mode,
                  const void        *alpha,
                  const cudnnTensorDescriptor_t biasDesc,
                  const void        *biasData,
                  const void        *beta,
                  const cudnnTensorDescriptor_t srcDestDesc,
                  void              *srcDestData )

```

This function adds the scaled values of one tensor to another tensor. The **mode** parameter can be used to select different ways of performing the scaled addition. The amount of data described by the **biasDesc** descriptor must match exactly the amount of data needed to perform the addition. Therefore, the following conditions must be met:

- ▶ Except for the **CUDNN\_ADD\_SAME\_C** mode, the dimensions **h**, **w** of the two tensors must match.
- ▶ In the case of **CUDNN\_ADD\_IMAGE** mode, the dimensions **n**, **c** of the bias tensor must be 1.
- ▶ In the case of **CUDNN\_ADD\_FEATURE\_MAP** mode, the dimension **n** of the bias tensor must be 1 and the dimension **c** of the two tensors must match.
- ▶ In the case of **CUDNN\_ADD\_FULL\_TENSOR** mode, the dimensions **n**, **c** of the two tensors must match.
- ▶ In the case of **CUDNN\_ADD\_SAME\_C** mode, the dimensions **n**, **w**, **h** of the bias tensor must be 1 and the dimension **c** of the two tensors must match.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
mode	input	Addition mode that describe how the addition is performed.
alpha	input	Scalar factor to be applied to every data element of the bias tensor before it is added to the output tensor.
biasDesc	input	Handle to a previously initialized tensor descriptor.
biasData	input	Pointer to data of the tensor described by the <b>biasDesc</b> descriptor.
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the operation Note that if <b>beta</b> is zero, the output is not read and can contain any uninitialized data (including Nan numbers)
srcDestDesc	input/ output	Handle to a previously initialized tensor descriptor.
srcDestData	input/ output	Pointer to data of the tensor described by the <b>srcDestDesc</b> descriptor.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The function executed successfully.
CUDNN_STATUS_BAD_PARAM	The dimensions <code>n, c, h, w</code> of the bias tensor refer to an amount of data that is incompatible with the <code>mode</code> parameter and the output tensor dimensions or the <code>dataType</code> of the two tensor descriptors are different.
CUDNN_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

## 4.16. cudnnSetTensor

```

cudnnStatus_t
cudnnSetTensor( cudnnHandle_t      handle,
                 const cudnnTensorDescriptor_t srcDestDesc,
                 void               *srcDestData,
                 const void         *value

```

This function sets all the elements of a tensor to a given value

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
srcDestDesc	input	Handle to a previously initialized tensor descriptor.
srcDestData	input/ output	Pointer to data of the tensor described by the <code>srcDestDesc</code> descriptor.
value	input	Pointer in Host memory to a value that all elements of the tensor will be set to.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The function launched successfully.
CUDNN_STATUS_BAD_PARAM	one of the provided pointers is nil
CUDNN_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

## 4.17. cudnnScaleTensor

```

cudnnStatus_t
cudnnScaleTensor( cudnnHandle_t      handle,
                  const cudnnTensorDescriptor_t srcDestDesc,
                  void               *srcDestData,
                  const void         *alpha

```

This function scale all the elements of a tensor by a give factor.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.

Param	In/out	Meaning
srcDestDesc	input	Handle to a previously initialized tensor descriptor.
srcDestData	input/ output	Pointer to data of the tensor described by the <code>srcDestDesc</code> descriptor.
alpha	input	Pointer in Host memory to a value that all elements of the tensor will be scaled with.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The function launched successfully.
<code>CUDNN_STATUS_BAD_PARAM</code>	one of the provided pointers is nil
<code>CUDNN_STATUS_EXECUTION_FAILED</code>	The function failed to launch on the GPU.

## 4.18. cudnnCreateFilterDescriptor

```

cudnnStatus_t cudnnCreateFilterDescriptor(cudnnFilterDescriptor_t *filterDesc)

```

This function creates a filter descriptor object by allocating the memory needed to hold its opaque structure,

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The object was created successfully.
<code>CUDNN_STATUS_ALLOC_FAILED</code>	The resources could not be allocated.

## 4.19. cudnnSetFilter4dDescriptor

```

cudnnStatus_t
cudnnSetFilter4dDescriptor( cudnnFilterDescriptor_t filterDesc,
                           cudnnDataType_t dataType,
                           int k,
                           int c,
                           int h,
                           int w )

```

This function initializes a previously created filter descriptor object into a 4D filter. Filters layout must be contiguous in memory.

Param	In/out	Meaning
filterDesc	input/ output	Handle to a previously created filter descriptor.
datatype	input	Data type.
k	input	Number of output feature maps.
c	input	Number of input feature maps.

Param	In/out	Meaning
h	input	Height of each filter.
w	input	Width of each filter.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was set successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the parameters <b>k</b> , <b>c</b> , <b>h</b> , <b>w</b> is negative or <b>dataType</b> has an invalid enumerant value.

## 4.20. cudnnGetFilter4dDescriptor

```

cudnnStatus_t
cudnnGetFilter4dDescriptor( cudnnFilterDescriptor_t filterDesc,
                           cudnnDataType_t *dataType,
                           int *k,
                           int *c,
                           int *h,
                           int *w )

```

This function queries the parameters of the previously initialized filter descriptor object.

Param	In/out	Meaning
filterDesc	input	Handle to a previously created filter descriptor.
datatype	output	Data type.
k	output	Number of output feature maps.
c	output	Number of input feature maps.
h	output	Height of each filter.
w	output	Width of each filter.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was set successfully.

## 4.21. cudnnSetFilterNdDescriptor

```

cudnnStatus_t
cudnnSetFilterNdDescriptor( cudnnFilterDescriptor_t filterDesc,
                           int nbDims,
                           int filterDimA[])

```

This function initializes a previously created filter descriptor object. Filters layout must be contiguous in memory.



Param	In/out	Meaning
filterDesc	input/ output	Handle to a previously created filter descriptor.
datatype	input	Data type.
nbDims	input	Dimension of the filter.
filterDimA	input	Array of dimension <code>nbDims</code> containing the size of the filter for each dimension.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was set successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the elements of the array <code>filterDimA</code> is negative or <code>dataType</code> has an invalid enumerant value.
CUDNN_STATUS_NOT_SUPPORTED	the parameter <code>nbDims</code> exceeds the maximum supported dimension.

## 4.22. cudnnGetFilterNdDescriptor

```

cudnnStatus_t
cudnnGetFilterNdDescriptor( const cudnnFilterDescriptor_t filterDesc,
                           int nbDimsRequested,
                           cudnnDataType_t *dataType,
                           int *nbDims,
                           int filterDimA[])

```

This function queries a previously initialized filter descriptor object.

Param	In/out	Meaning
filterDesc	input	Handle to a previously initialized filter descriptor.
nbDimsRequested	input	Dimension of the expected filter descriptor. It is also the minimum size of the arrays <code>filterDimA</code> in order to be able to hold the results
datatype	input	Data type.
nbDims	input	Actual dimension of the filter.
filterDimA	input	Array of dimension of at least <code>nbDimsRequested</code> that will be filled with the filter parameters from the provided filter descriptor.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was set successfully.
CUDNN_STATUS_BAD_PARAM	The parameter <code>nbDimsRequested</code> is negative.

## 4.23. cudnnDestroyFilterDescriptor

```
cudaStatus_t cudnnDestroyFilterDescriptor(cudaFilterDescriptor_t filterDesc)
```

This function destroys a previously created Tensor4D descriptor object.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was destroyed successfully.

## 4.24. cudnnCreateConvolutionDescriptor

```
cudaStatus_t cudnnCreateConvolutionDescriptor(cudaConvolutionDescriptor_t *convDesc)
```

This function creates a convolution descriptor object by allocating the memory needed to hold its opaque structure,

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was created successfully.
CUDNN_STATUS_ALLOC_FAILED	The resources could not be allocated.

## 4.25. cudnnSetConvolution2dDescriptor

```
cudaStatus_t
cudnnSetConvolution2dDescriptor( cudaConvolutionDescriptor_t convDesc,
                                int pad_h,
                                int pad_w,
                                int u,
                                int v,
                                int upscalex,
                                int upscaley,
                                cudaConvolutionMode_t mode )
```

This function initializes a previously created convolution descriptor object into a 2D correlation. This function assumes that the tensor and filter descriptors corresponds to the forward convolution path and checks if their settings are valid. That same convolution descriptor can be reused in the backward path provided it corresponds to the same layer.

Param	In/out	Meaning
convDesc	input/output	Handle to a previously created convolution descriptor.
pad_h	input	zero-padding height: number of rows of zeros implicitly concatenated onto the top and onto the bottom of input images.
pad_w	input	zero-padding width: number of columns of zeros implicitly concatenated onto the left and onto the right of input images.
u	input	Vertical filter stride.

Param	In/out	Meaning
v	input	Horizontal filter stride.
upscalex	input	Upscale the input in x-direction.
upscaley	input	Upscale the input in y-direction.
mode	input	Selects between <code>CUDNN_CONVOLUTION</code> and <code>CUDNN_CROSS_CORRELATION</code> .

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The object was set successfully.
<code>CUDNN_STATUS_BAD_PARAM</code>	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>► One of the parameters <code>u</code>, <code>v</code> is negative.</li> <li>► The parameter <code>mode</code> has an invalid enumerant value.</li> </ul>
<code>CUDNN_STATUS_NOT_SUPPORTED</code>	The parameter <code>upscalex</code> or <code>upscaley</code> is not 1.

## 4.26. cudnnGetConvolution2dDescriptor

```

cudnnStatus_t
cudnnGetConvolution2dDescriptor( const cudnnConvolutionDescriptor_t convDesc,
                                int* pad_h,
                                int* pad_w,
                                int* u,
                                int* v,
                                int* upscalex,
                                int* upscaley,
                                cudnnConvolutionMode_t *mode )

```

This function queries a previously initialized 2D convolution descriptor object.

Param	In/out	Meaning
convDesc	input/ output	Handle to a previously created convolution descriptor.
pad_h	output	zero-padding height: number of rows of zeros implicitly concatenated onto the top and onto the bottom of input images.
pad_w	output	zero-padding width: number of columns of zeros implicitly concatenated onto the left and onto the right of input images.
u	output	Vertical filter stride.
v	output	Horizontal filter stride.
upscalex	output	Upscale the input in x-direction.
upscaley	output	Upscale the input in y-direction.
mode	output	convolution mode.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was set successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>▶ One of the parameters <code>u</code>, <code>v</code> is negative.</li> <li>▶ The parameter <code>mode</code> has an invalid enumerant value.</li> </ul>
CUDNN_STATUS_NOT_SUPPORTED	The parameter <code>upscalex</code> or <code>upscaley</code> is not 1.

## 4.27. cudnnGetConvolution2dForwardOutputDim

```

cudnnStatus_t
cudnnGetConvolution2dForwardOutputDim( const cudnnConvolutionDescriptor_t
    convDesc,
                                        const cudnnTensorDescriptor_t
    inputTensorDesc,
                                        const cudnnFilterDescriptor_t filterDesc,
    int *n,
    int *c,
    int *h,
    int *w )

```

This function returns the dimensions of the resulting 4D tensor of a 2D convolution, given the convolution descriptor, the input tensor descriptor and the filter descriptor. This function can help to setup the output tensor and allocate the proper amount of memory prior to launch the actual convolution.

Each dimension **h** and **w** of the output images is computed as followed:

```
outputDim = 1 + (inputDim + 2*pad - filterDim)/convolutionStride;
```

Param	In/out	Meaning
convDesc	input	Handle to a previously created convolution descriptor.
inputTensorDesc	input	Handle to a previously initialized tensor descriptor.
filterDesc	input	Handle to a previously initialized filter descriptor.
n	output	Number of output images.
c	output	Number of output feature maps per image.
h	output	Height of each output feature map.
w	output	Width of each output feature map.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_BAD_PARAM	One or more of the descriptors has not been created correctly or there is a mismatch between

Return Value	Meaning
	the feature maps of <code>inputTensorDesc</code> and <code>filterDesc</code> .
<code>CUDNN_STATUS_SUCCESS</code>	The object was set successfully.

## 4.28. cudnnSetConvolutionNdDescriptor

```

cudnnStatus_t
cudnnSetConvolutionNdDescriptor( cudnnConvolutionDescriptor_t convDesc,
                                int arrayLength,
                                int padA[],

                                int filterStrideA[],
                                int upscaleA[],
                                cudnnConvolutionMode_t mode )

```

This function initializes a previously created generic convolution descriptor object into a n-D correlation. That same convolution descriptor can be reused in the backward path provided it corresponds to the same layer.

Param	In/out	Meaning
<code>convDesc</code>	input/output	Handle to a previously created convolution descriptor.
<code>arrayLength</code>	input	Dimension of the convolution.
<code>padA</code>	input	Array of dimension <code>arrayLength</code> containing the zero-padding size for each dimension. For every dimension, the padding represents the number of extra zeros implicitly concatenated at the start and at the end of every element of that dimension .
<code>filterStrideA</code>	input	Array of dimension <code>arrayLength</code> containing the filter stride for each dimension. For every dimension, the filter stride represents the number of elements to slide to reach the next start of the filtering window of the next point.
<code>upscaleA</code>	input	Array of dimension <code>arrayLength</code> containing the upscale factor for each dimension.
<code>mode</code>	input	Selects between <code>CUDNN_CONVOLUTION</code> and <code>CUDNN_CROSS_CORRELATION</code> .

## 4.29. cudnnGetConvolutionNdDescriptor

```

cudnnStatus_t
cudnnGetConvolutionNdDescriptor( const cudnnConvolutionDescriptor_t convDesc,
                                int arrayLengthRequested,
                                int *arrayLength,
                                int padA[],

                                int filterStrideA[],
                                int upscaleA[],
                                cudnnConvolutionMode_t *mode )

```

This function queries a previously initialized convolution descriptor object.

Param	In/out	Meaning
convDesc	input/output	Handle to a previously created convolution descriptor.
arrayLengthRequested	input	Dimension of the expected convolution descriptor. It is also the minimum size of the arrays <code>padA</code> , <code>filterStrideA</code> and <code>upscaleA</code> in order to be able to hold the results
arrayLength	output	actual dimension of the convolution descriptor.
padA	output	Array of dimension of at least <code>arrayLengthRequested</code> that will be filled with the padding parameters from the provided convolution descriptor.
filterStrideA	output	Array of dimension of at least <code>arrayLengthRequested</code> that will be filled with the filter stride from the provided convolution descriptor.
upscaleA	output	Array of dimension at least <code>arrayLengthRequested</code> that will be filled with the upscaling parameters from the provided convolution descriptor.
mode	output	convolution mode of the provided descriptor.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The query was successfully.
CUDNN_STATUS_BAD_PARAM	The <code>arrayLengthRequest</code> is negative
CUDNN_STATUS_NOT_SUPPORTED	The <code>arrayLengthRequest</code> is greater than the maximum supported dimension

## 4.30. cudnnGetConvolutionNdForwardOutputDim

```

cudnnStatus_t
cudnnGetConvolutionNdForwardOutputDim( const cudnnConvolutionDescriptor_t
    convDesc,
                                        const cudnnTensorDescriptor_t
    inputTensorDesc,
                                        const cudnnFilterDescriptor_t filterDesc,
                                        int nbDims,
                                        int tensorOutputDimA[] )

```

This function returns the dimensions of the resulting n-D tensor of a **nbDims-2-D** convolution, given the convolution descriptor, the input tensor descriptor and the filter descriptor. This function can help to setup the output tensor and allocate the proper amount of memory prior to launch the actual convolution.

Each dimension of the **(nbDims-2) -D** images of the output tensor is computed as followed:

```
outputDim = 1 + (inputDim + 2*pad - filterDim)/convolutionStride;
```

Param	In/out	Meaning
convDesc	input	Handle to a previously created convolution descriptor.
inputTensorDesc	input	Handle to a previously initialized tensor descriptor.
filterDesc	input	Handle to a previously initialized filter descriptor.
nbDims	input	Dimension of the output tensor
tensorOutputDims	output	Array of dimensions <code>nbDims</code> that contains on exit of this routine the sizes of the output tensor

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
<code>CUDNN_STATUS_BAD_PARAM</code>	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>▶ The dimension of the convolution descriptor is different from <code>nbDims-2</code>.</li> <li>▶ The dimension of the input tensor descriptor is different from <code>nbDims</code>.</li> <li>▶ The dimension of the filter descriptor is different from <code>nbDims</code>.</li> </ul>
<code>CUDNN_STATUS_SUCCESS</code>	The routine exits successfully.

## 4.31. cudnnDestroyFilterDescriptor

```

cudnnStatus_t cudnnDestroyConvolutionDescriptor(cudnnConvolutionDescriptor_t
convDesc)

```

This function destroys a previously created convolution descriptor object.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The object was destroyed successfully.

## 4.32. cudnnGetConvolutionForwardAlgorithm

```

cudnnStatus_t
cudnnGetConvolutionForwardAlgorithm( cudnnHandle_t          handle,
                                     const cudnnTensorDescriptor_t srcDesc,
                                     const cudnnFilterDescriptor_t
filterDesc,
                                     const cudnnConvolutionDescriptor_t
convDesc,
                                     const cudnnTensorDescriptor_t
destDesc,
                                     cudnnConvolutionFwdPreference_t
preference,
                                     size_t
memoryLimitInbytes,
                                     cudnnConvolutionFwdAlgo_t    *algo
)

```

This function advised the best algorithm to choose for the forward convolution depending on the criteria expressed in the `cudaConvolutionFwdPreference_t` enumerant.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
srcDesc	input	Handle to the previously initialized input tensor descriptor.
filterDesc	input	Handle to a previously initialized filter descriptor.
convDesc	input	Previously initialized convolution descriptor.
destDesc	input	Handle to the previously initialized output tensor descriptor.
preference	input	Enumerant to express the preference criteria in terms of memory requirement and speed.
memoryLimitInBytes	input	It is used when enumerant <code>preference</code> is set to <code>CUDNN_CONVOLUTION_FWD_SPECIFY_WORKSPACE_LIMIT</code> to specify the maximum amount of GPU memory the user is willing to use as a workspace
algo	output	Enumerant that specifies which convolution algorithm should be used to compute the results according to the specified preference

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The query was successful.
<code>CUDNN_STATUS_BAD_PARAM</code>	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>► The numbers of feature maps of the input tensor and output tensor differ.</li> <li>► The <code>dataType</code> of the two tensor descriptors or the filter are different.</li> </ul>

## 4.33. cudnnGetConvolutionForwardWorkspaceSize

```

cudaStatus_t
cudnnGetConvolutionForwardWorkspaceSize( cudaHandle_t    handle,
                                         const  cudaTensorDescriptor_t
srcDesc,                                         const  cudaFilterDescriptor_t
                                         const  cudaConvolutionDescriptor_t
filterDesc,                                     const  cudaTensor4dDescriptor_t
                                         const  cudaTensor4dDescriptor_t
convDesc,                                     cudaConvolutionFwdAlgo_t
destDesc,
algo,
size_t
*sizeInBytes
)

```

This function returns the amount of GPU memory workspace the user needs to allocate to be able to call `cudaConvolutionForward` with the specified



algorithm. The workspace allocated will then be passed to the routine **cudaConvolutionForward**. The specified algorithm can be the result of the call to **cudaGetConvolutionForwardAlgorithm** or can be chosen arbitrarily by the user. Note that not every algorithm is available for every configuration of the input tensor and/or every configuration of the convolution descriptor.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
srcDesc	input	Handle to the previously initialized input tensor descriptor.
filterDesc	input	Handle to a previously initialized filter descriptor.
convDesc	input	Previously initialized convolution descriptor.
destDesc	input	Handle to the previously initialized output tensor descriptor.
algo	input	Enumerant that specifies the chosen convolution algorithm
sizeInBytes	output	Amount of GPU memory needed as workspace to be able to execute a forward convolution with the specified <b>algo</b>

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The query was successful.
CUDNN_STATUS_BAD_PARAM	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>► The numbers of feature maps of the input tensor and output tensor differ.</li> <li>► The <b>dataType</b> of the two tensor descriptors or the filter are different.</li> </ul>
CUDNN_STATUS_NOT_SUPPORTED	The combination of the tensor descriptors, filter descriptor and convolution descriptor is not supported for the specified algorithm.

## 4.34. cudaConvolutionForward

```

cudaStatus_t
cudaConvolutionForward( cudaHandle_t          handle,
                        const void            *alpha,
                        const cudnnTensorDescriptor_t srcDesc,
                        const void            *srcData,
                        const cudnnFilterDescriptor_t filterDesc,
                        const void            *filterData,
                        const cudnnConvolutionDescriptor_t convDesc,
                        cudnnConvolutionFwdAlgo_t algo,
                        void                  *workspace,
                        size_t                size_t,
                        workspaceSizeInBytes,
                        const void            *beta,
                        const cudnnTensorDescriptor_t destDesc,
                        void                  *destData )

```

This function executes convolutions or cross-correlations over **src** using the specified **filters**, returning results in **dest**. Scaling factors **alpha** and **beta** can be used to scale the input tensor and the output tensor respectively.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
alpha	input	Scaling factor with which every element of the input tensor is multiplied.
srcDesc	input	Handle to a previously initialized tensor descriptor.
srcData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDesc</b> .
filterDesc	input	Handle to a previously initialized filter descriptor.
filterData	input	Data pointer to GPU memory associated with the filter descriptor <b>filterDesc</b> .
convDesc	input	Previously initialized convolution descriptor.
algo	input	Enumerant that specifies which convolution algorithm should be used to compute the results
workSpace	input	Data pointer to GPU memory to a workspace needed to able to execute the specified algorithm. If no workspace is needed for a particular algorithm, that pointer can be nil
workSpaceSizeInBytes	input	Specifies the size in bytes of the provided <b>workSpace</b>
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the convolution. Note that if <b>beta</b> is zero, the output is not read and can contain any uninitialized data (including Nan numbers)
destDesc	input	Handle to a previously initialized tensor descriptor.
destData	input/ output	Data pointer to GPU memory associated with the tensor descriptor <b>destDesc</b> that carries the result of the convolution.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The operation was launched successfully.
CUDNN_STATUS_MAPPING_ERROR	An error occurred during the texture binding of the filter data.
CUDNN_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

## 4.35. cudnnConvolutionBackwardBias

```

cudnnStatus_t
cudnnConvolutionBackwardBias( cudnnHandle_t      handle,
                              const void*        *alpha,
                              const cudnnTensorDescriptor_t srcDesc,
                              const void*        *srcData,
                              const void*        *beta,
                              const cudnnTensorDescriptor_t destDesc,
                              void*              *destData
                              )

```

This function computes the convolution gradient with respect to the bias, which is the sum of every element belonging to the same feature map across all of the images of the input tensor. Therefore, the number of elements produced is equal to the number of features maps of the input tensor.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
alpha	input	Scaling factor with which every element of the input tensor is multiplied.
srcDesc	input	Handle to the previously initialized input tensor descriptor.
srcData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDesc</b> .
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the convolution gradient. Note that if <b>beta</b> is zero, the output is not read and can contain any uninitialized data (including Nan numbers)
destDesc	input	Handle to the previously initialized output tensor descriptor.
destData	output	Data pointer to GPU memory associated with the output tensor descriptor <b>destDesc</b> .

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The operation was launched successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>One of the parameters <b>n, h, w</b> of the output tensor is not 1.</li> <li>The numbers of feature maps of the input tensor and output tensor differ.</li> <li>The <b>dataType</b> of the two tensor descriptors are different.</li> </ul>

## 4.36. cudnnConvolutionBackwardFilter

```

cudnnStatus_t
cudnnConvolutionBackwardFilter( cudnnHandle_t      handle,
                                const void         *alpha,
                                const cudnnTensorDescriptor_t srcDesc,
                                const void         *srcData,
                                const cudnnTensorDescriptor_t diffDesc,
                                const void         *diffData,
                                const cudnnConvolutionDescriptor_t convDesc,
                                const void         *beta,

                                const cudnnFilterDescriptor_t gradDesc,
                                void              *gradData )

```

This function computes the convolution gradient with respect to the filter coefficients.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
alpha	input	Scaling factor with which every element of the input tensor is multiplied.
srcDesc	input	Handle to a previously initialized tensor descriptor.
srcData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDesc</b> .
diffDesc	input	Handle to the previously initialized input differential tensor descriptor.
diffData	input	Data pointer to GPU memory associated with the input differential tensor descriptor <b>diffDesc</b> .
convDesc	input	Previously initialized convolution descriptor.
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the convolution gradient. Note that if <b>beta</b> is zero, the output is not read and can contain any uninitialized data (including Nan numbers)
gradDesc	input	Handle to a previously initialized filter descriptor.
gradData	input/ output	Data pointer to GPU memory associated with the filter descriptor <b>gradDesc</b> that carries the result.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The operation was launched successfully.
CUDNN_STATUS_NOT_SUPPORTED	The requested operation is not currently supported in cuDNN. The descriptor <b>diffDesc</b> is likely not in NCHW format.
CUDNN_STATUS_MAPPING_ERROR	An error occurs during the texture binding of the filter data.
CUDNN_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

## 4.37. cudnnConvolutionBackwardData

```

cudnnStatus_t
cudnnConvolutionBackwardData( cudnnHandle_t      handle,
                              const void         *alpha,
                              const cudnnFilterDescriptor_t filterDesc,
                              const void         *filterData,
                              const cudnnTensorDescriptor_t diffDesc,
                              const void         *diffData,
                              const cudnnConvolutionDescriptor_t convDesc,
                              const void         *beta,
                              const cudnnTensorDescriptor_t gradDesc,
                              void               *gradData
                              );

```

This function computes the convolution gradient with respect to the output tensor.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
alpha	input	Scaling factor with which every element of the input tensor is multiplied.
filterDesc	input	Handle to a previously initialized filter descriptor.
filterData	input	Data pointer to GPU memory associated with the filter descriptor <b>filterDesc</b> .
diffDesc	input	Handle to the previously initialized input differential tensor descriptor.
diffData	input	Data pointer to GPU memory associated with the input differential tensor descriptor <b>diffDesc</b> .
convDesc	input	Previously initialized convolution descriptor.
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the convolution gradient. Note that if <b>beta</b> is zero, the output is not read and can contain any uninitialized data (including Nan numbers)
gradDesc	input	Handle to the previously initialized output tensor descriptor.
gradData	input/ output	Data pointer to GPU memory associated with the output tensor descriptor <b>gradDesc</b> that carries the result.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The operation was launched successfully.
CUDNN_STATUS_NOT_SUPPORTED	The requested operation is not currently supported in cuDNN. The descriptor <b>diffDesc</b> is likely not in NCHW format.
CUDNN_STATUS_MAPPING_ERROR	An error occurs during the texture binding of the filter data or the input differential tensor data

Return Value	Meaning
CUDNN_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

## 4.38. cudnnSoftmaxForward

```

cudnnStatus_t
cudnnSoftmaxForward( cudnnHandle_t          handle,
                     cudnnSoftmaxAlgorithm_t algorithm,
                     cudnnSoftmaxMode_t      mode,
                     const void              *alpha,
                     const cudnnTensorDescriptor_t srcDesc,
                     const void              *srcData,
                     const void              *beta,
                     const cudnnTensorDescriptor_t destDesc,
                     void                    *destData )

```

This routine computes the softmax function.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
algorithm	input	Enumerant to specify the softmax algorithm.
mode	input	Enumerant to specify the softmax mode.
alpha	input	Scaling factor with which every element of the input tensor is multiplied.
srcDesc	input	Handle to the previously initialized input tensor descriptor.
srcData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDesc</b> .
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the softmax function. Note that if <b>beta</b> is zero, the output is not read and can contain any uninitialized data (including Nan numbers)
destDesc	input	Handle to the previously initialized output tensor descriptor.
destData	output	Data pointer to GPU memory associated with the output tensor descriptor <b>destDesc</b> .

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The function launched successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>▶ The dimensions <b>n, c, h, w</b> of the input tensor and output tensors differ.</li> <li>▶ The <b>datatype</b> of the input tensor and output tensors differ.</li> <li>▶ The parameters <b>algorithm</b> or <b>mode</b> have an invalid enumerant value.</li> </ul>

Return Value	Meaning
CUDNN_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

## 4.39. cudnnSoftmaxBackward

```

cudnnStatus_t
cudnnSoftmaxBackward( cudnnHandle_t          handle,
                      cudnnSoftmaxAlgorithm_t algorithm,
                      cudnnSoftmaxMode_t      mode,
                      const void              *alpha,
                      const cudnnTensorDescriptor_t srcDesc,
                      const void              *srcData,
                      const cudnnTensorDescriptor_t srcDiffDesc,
                      const void              *srcDiffData,
                      const void              *beta,
                      const cudnnTensorDescriptor_t destDiffDesc,
                      void                    *destDiffData )

```

This routine computes the gradient of the softmax function.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
algorithm	input	Enumerant to specify the softmax algorithm.
mode	input	Enumerant to specify the softmax mode.
alpha	input	Scaling factor with which every element of the input tensors is multiplied.
srcDesc	input	Handle to the previously initialized input tensor descriptor.
srcData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDesc</b> .
srcDiffDesc	input	Handle to the previously initialized input differential tensor descriptor.
srcDiffData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDiffData</b> .
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the softmax gradient. Note that if <b>beta</b> is zero, the output is not read and can contain any uninitialized data (including Nan numbers)
destDiffDesc	input	Handle to the previously initialized output differential tensor descriptor.
destDiffData	output	Data pointer to GPU memory associated with the output tensor descriptor <b>destDiffDesc</b> .

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The function launched successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the following conditions are met:

Return Value	Meaning
	<ul style="list-style-type: none"> <li>▶ The dimensions <i>n, c, h, w</i> of the <i>srcDesc</i>, <i>srcDiffDesc</i> and <i>destDiffDesc</i> tensors differ.</li> <li>▶ The strides <i>nStride</i>, <i>cStride</i>, <i>hStride</i>, <i>wStride</i> of the <i>srcDesc</i> and <i>srcDiffDesc</i> tensors differ.</li> <li>▶ The <i>datatype</i> of the three tensors differs.</li> </ul>
<b>CUDNN_STATUS_EXECUTION_FAILED</b>	The function failed to launch on the GPU.

## 4.40. cudnnCreatePoolingDescriptor

```

cudnnStatus_t cudnnCreatePoolingDescriptor( cudnnPoolingDescriptor_t*
poolingDesc )

```

This function creates a pooling descriptor object by allocating the memory needed to hold its opaque structure,

Return Value	Meaning
<b>CUDNN_STATUS_SUCCESS</b>	The object was created successfully.
<b>CUDNN_STATUS_ALLOC_FAILED</b>	The resources could not be allocated.

## 4.41. cudnnSetPooling2dDescriptor

```

cudnnStatus_t
cudnnSetPooling2dDescriptor( cudnnPoolingDescriptor_t poolingDesc,
                             cudnnPoolingMode_t mode,
                             int windowHeight,
                             int windowWidth,
                             int verticalPadding,
                             int horizontalPadding,
                             int verticalStride,
                             int horizontalStride )

```

This function initializes a previously created generic pooling descriptor object into a 2D description.

Param	In/out	Meaning
<i>poolingDesc</i>	input/output	Handle to a previously created pooling descriptor.
<i>mode</i>	input	Enumerant to specify the pooling mode.
<i>windowHeight</i>	input	Height of the pooling window.
<i>windowWidth</i>	input	Width of the pooling window.
<i>verticalPadding</i>	input	Size of vertical padding.
<i>horizontalPadding</i>	input	Size of horizontal padding
<i>verticalStride</i>	input	Pooling vertical stride.



Param	In/out	Meaning
horizontalStride	input	Pooling horizontal stride.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was set successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the parameters <code>windowHeight</code> , <code>windowWidth</code> , <code>verticalStride</code> , <code>horizontalStride</code> is negative or <code>mode</code> has an invalid enumerant value.

## 4.42. cudnnGetPooling2dDescriptor

```

cudnnStatus_t
cudnnGetPooling2dDescriptor( const cudnnPoolingDescriptor_t poolingDesc,
                             cudnnPoolingMode_t *mode,
                             int *windowHeight,
                             int *windowWidth,
                             int *verticalPadding,
                             int *horizontalPadding,
                             int *verticalStride,
                             int *horizontalStride )

```

This function queries a previously created 2D pooling descriptor object.

Param	In/out	Meaning
poolingDesc	input	Handle to a previously created pooling descriptor.
mode	output	Enumerant to specify the pooling mode.
windowHeight	output	Height of the pooling window.
windowWidth	output	Width of the pooling window.
verticalPadding	output	Size of vertical padding.
horizontalPadding	output	Size of horizontal padding.
verticalStride	output	Pooling vertical stride.
horizontalStride	output	Pooling horizontal stride.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was set successfully.

## 4.43. cudnnSetPoolingNdDescriptor

```

cudnnStatus_t
cudnnSetPoolingNdDescriptor( cudnnPoolingDescriptor_t poolingDesc,
                             cudnnPoolingMode_t mode,
                             int nbDims,
                             int windowDimA[],
                             int paddingA[],
                             int strideA[] )

```

This function initializes a previously created generic pooling descriptor object.

Param	In/out	Meaning
poolingDesc	input/ output	Handle to a previously created pooling descriptor.
mode	input	Enumerant to specify the pooling mode.
nbDims	input	Dimension of the pooling operation.
windowDimA	output	Array of dimension <b>nbDims</b> containing the window size for each dimension.
paddingA	output	Array of dimension <b>nbDims</b> containing the padding size for each dimension.
strideA	output	Array of dimension <b>nbDims</b> containing the striding size for each dimension.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The object was set successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the elements of the arrays <b>windowDimA</b> , <b>paddingA</b> or <b>strideA</b> is negative or <b>mode</b> has an invalid enumerant value.

## 4.44. cudnnGetPoolingNdDescriptor

```

cudnnStatus_t
cudnnGetPoolingNdDescriptor( const cudnnPoolingDescriptor_t poolingDesc,
                             int nbDimsRequested,
                             cudnnPoolingMode_t *mode,
                             int *nbDims,
                             int windowDimA[],
                             int paddingA[],
                             int strideA[] )

```

This function queries a previously initialized generic pooling descriptor object.

Param	In/out	Meaning
poolingDesc	input	Handle to a previously created pooling descriptor.

Param	In/out	Meaning
nbDimsRequested	input	Dimension of the expected pooling descriptor. It is also the minimum size of the arrays <code>windowDimA</code> , <code>paddingA</code> and <code>strideA</code> in order to be able to hold the results
mode	output	Enumerant to specify the pooling mode.
nbDims	output	Actual dimension of the pooling descriptor.
windowDimA	output	Array of dimension of at least <code>nbDimsRequested</code> that will be filled with the window parameters from the provided pooling descriptor.
paddingA	output	Array of dimension of at least <code>nbDimsRequested</code> that will be filled with the padding parameters from the provided pooling descriptor.
strideA	output	Array of dimension at least <code>nbDimsRequested</code> that will be filled with the stride parameters from the provided pooling descriptor.

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The object was queried successfully.
<code>CUDNN_STATUS_BAD_PARAM</code>	The parameter <code>nbDimsRequested</code> is negative.

## 4.45. cudnnDestroyPoolingDescriptor

```

cudnnStatus_t cudnnDestroyPoolingDescriptor( cudnnPoolingDescriptor_t
poolingDesc )

```

This function destroys a previously created pooling descriptor object.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The object was destroyed successfully.

## 4.46. cudnnPoolingForward

```

cudnnStatus_t
cudnnPoolingForward( cudnnHandle_t          handle,
                    const cudnnPoolingDescriptor_t poolingDesc,
                    const void              *alpha,
                    const cudnnTensorDescriptor_t srcDesc,
                    const void              *srcData,
                    const void              *beta,
                    const cudnnTensorDescriptor_t destDesc,
                    void                    *destData )

```

This function computes pooling of input values (i.e., the maximum or average of several adjacent values) to produce an output with smaller height and/or width.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.

Param	In/out	Meaning
poolingDesc	input	Handle to a previously initialized pooling descriptor.
alpha	input	Scaling factor with which every element of the input tensor is multiplied.
srcDesc	input	Handle to the previously initialized input tensor descriptor.
srcData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDesc</b> .
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the pooling. Note that if <b>beta</b> is zero, the output is not read and can contain any uninitialized data (including Nan numbers)
destDesc	input	Handle to the previously initialized output tensor descriptor.
destData	output	Data pointer to GPU memory associated with the output tensor descriptor <b>destDesc</b> .

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The function launched successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>► The dimensions <b>n, c</b> of the input tensor and output tensors differ.</li> <li>► The <b>datatype</b> of the input tensor and output tensors differs.</li> </ul>
CUDNN_STATUS_NOT_SUPPORTED	The <b>wStride</b> of input tensor or output tensor is not 1.
CUDNN_STATUS_EXECUTION_FAILED	The function failed to launch on the GPU.

## 4.47. cudnnPoolingBackward

```

cudnnStatus_t
cudnnPoolingBackward( cudnnHandle_t handle,
    const cudnnPoolingDescriptor_t poolingDesc,
    const void *alpha,
    const cudnnTensorDescriptor_t srcDesc,
    const void *srcData,
    const cudnnTensorDescriptor_t srcDiffDesc,
    const void *srcDiffData,
    const cudnnTensorDescriptor_t destDesc,
    const void *destData,
    const void *beta,
    const cudnnTensorDescriptor_t destDiffDesc,
    void *destDiffData )

```

This function computes the gradient of a pooling operation.

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.

Param	In/out	Meaning
poolingDesc	input	Handle to the previously initialized pooling descriptor.
alpha	input	Scaling factor with which every element of the input tensors is multiplied.
srcDesc	input	Handle to the previously initialized input tensor descriptor.
srcData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDesc</b> .
srcDiffDesc	input	Handle to the previously initialized input differential tensor descriptor.
srcDiffData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDiffData</b> .
destDesc	input	Handle to the previously initialized output tensor descriptor.
destData	input	Data pointer to GPU memory associated with the output tensor descriptor <b>destDesc</b> .
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the pooling gradient. Note that if <b>beta</b> is zero, the output is not read and can contain any uninitialized data (including Nan numbers)
destDiffDesc	input	Handle to the previously initialized output differential tensor descriptor.
destDiffData	output	Data pointer to GPU memory associated with the output tensor descriptor <b>destDiffDesc</b> .

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
<b>CUDNN_STATUS_SUCCESS</b>	The function launched successfully.
<b>CUDNN_STATUS_BAD_PARAM</b>	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>► The dimensions <b>n, c, h, w</b> of the <b>srcDesc</b> and <b>srcDiffDesc</b> tensors differ.</li> <li>► The strides <b>nStride, cStride, hStride, wStride</b> of the <b>srcDesc</b> and <b>srcDiffDesc</b> tensors differ.</li> <li>► The dimensions <b>n, c, h, w</b> of the <b>destDesc</b> and <b>destDiffDesc</b> tensors differ.</li> <li>► The strides <b>nStride, cStride, hStride, wStride</b> of the <b>destDesc</b> and <b>destDiffDesc</b> tensors differ.</li> <li>► The <b>datatype</b> of the four tensors differ.</li> </ul>
<b>CUDNN_STATUS_NOT_SUPPORTED</b>	The <b>wStride</b> of input tensor or output tensor is not 1.
<b>CUDNN_STATUS_EXECUTION_FAILED</b>	The function failed to launch on the GPU.

## 4.48. cudnnActivationForward

```

cudnnStatus_t
cudnnActivationForward( cudnnHandle_t      handle,
                       cudnnActivationMode_t mode,
                       const void          *alpha,
                       const cudnnTensorDescriptor_t srcDesc,
                       const void          *srcData,
                       const void          *beta,
                       const cudnnTensorDescriptor_t destDesc,
                       void                *destData )

```

This routine applies a specified neuron activation function element-wise over each input value.



In-place operation is allowed for this routine; i.e., **srcData** and **destData** pointers may be equal. However, this requires **srcDesc** and **destDesc** descriptors to be identical (particularly, the strides of the input and output must match for in-place operation to be allowed).

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
mode	input	Enumerant to specify the activation mode.
alpha	input	Scaling factor with which every element of the input tensor is multiplied.
srcDesc	input	Handle to the previously initialized input tensor descriptor.
srcData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDesc</b> .
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the activation. Note that if <b>beta</b> is zero, the output is not read and can contain any uninitialized data (including Nan numbers).
destDesc	input	Handle to the previously initialized output tensor descriptor.
destData	output	Data pointer to GPU memory associated with the output tensor descriptor <b>destDesc</b> .

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
CUDNN_STATUS_SUCCESS	The function launched successfully.
CUDNN_STATUS_BAD_PARAM	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>▶ The parameter <b>mode</b> has an invalid enumerant value.</li> <li>▶ The dimensions <b>n, c, h, w</b> of the input tensor and output tensors differ.</li> <li>▶ The <b>datatype</b> of the input tensor and output tensors differs.</li> </ul>

Return Value	Meaning
	<ul style="list-style-type: none"> <li>The strides <b>nStride</b>, <b>cStride</b>, <b>hStride</b>, <b>wStride</b> of the input tensor and output tensors differ and in-place operation is used (i.e., <b>srcData</b> and <b>destData</b> pointers are equal).</li> </ul>
<b>CUDNN_STATUS_EXECUTION_FAILED</b>	The function failed to launch on the GPU.

## 4.49. cudnnActivationBackward

```

cudnnStatus_t
cudnnActivationBackward( cudnnHandle_t          handle,
                        cudnnActivationMode_t    mode,
                        const void              *alpha,
                        const cudnnTensorDescriptor_t srcDesc,
                        const void              *srcData,
                        const cudnnTensorDescriptor_t srcDiffDesc,
                        const void              *srcDiffData,
                        const cudnnTensorDescriptor_t destDesc,
                        const void              *destData,
                        const void              *beta,
                        const cudnnTensorDescriptor_t destDiffDesc,
                        void                    *destDiffData )

```

This routine computes the gradient of a neuron activation function.



In-place operation is allowed for this routine; i.e. **srcDiffData** and **destDiffData** pointers may be equal. However, this requires the corresponding tensor descriptors to be identical (particularly, the strides of the input and output must match for in-place operation to be allowed).

Param	In/out	Meaning
handle	input	Handle to a previously created cuDNN context.
mode	input	Enumerant to specify the activation mode.
alpha	input	Scaling factor with which every element of the input tensor is multiplied.
srcDesc	input	Handle to the previously initialized input tensor descriptor.
srcData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDesc</b> .
srcDiffDesc	input	Handle to the previously initialized input differential tensor descriptor.
srcDiffData	input	Data pointer to GPU memory associated with the tensor descriptor <b>srcDiffData</b> .
destDesc	input	Handle to the previously initialized output tensor descriptor.
destData	input	Data pointer to GPU memory associated with the output tensor descriptor <b>destDesc</b> .
beta	input	Scaling factor which is applied on every element of the output tensor prior to adding the result of the activation gradient Note that if <b>beta</b> is zero,

Param	In/out	Meaning
		the output is not read and can contain any uninitialized data (including Nan numbers)
destDiffDesc	input	Handle to the previously initialized output differential tensor descriptor.
destDiffData	output	Data pointer to GPU memory associated with the output tensor descriptor <code>destDiffDesc</code> .

The possible error values returned by this function and their meanings are listed below.

Return Value	Meaning
<code>CUDNN_STATUS_SUCCESS</code>	The function launched successfully.
<code>CUDNN_STATUS_BAD_PARAM</code>	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>▶ The parameter <code>mode</code> has an invalid enumerant value.</li> <li>▶ The dimensions <code>n</code>, <code>c</code>, <code>h</code>, <code>w</code> of the input tensor and output tensors differ.</li> <li>▶ The <code>datatype</code> of the input tensor and output tensors differs.</li> <li>▶ The strides <code>nStride</code>, <code>cStride</code>, <code>hStride</code>, <code>wStride</code> of the input differential tensor and output differential tensors differ and in-place operation is used.</li> </ul>
<code>CUDNN_STATUS_NOT_SUPPORTED</code>	At least one of the following conditions are met: <ul style="list-style-type: none"> <li>▶ The strides <code>nStride</code>, <code>cStride</code>, <code>hStride</code>, <code>wStride</code> of the input tensor and the input differential tensor differ.</li> <li>▶ The strides <code>nStride</code>, <code>cStride</code>, <code>hStride</code>, <code>wStride</code> of the output tensor and the output differential tensor differ.</li> </ul>
<code>CUDNN_STATUS_EXECUTION_FAILED</code>	The function failed to launch on the GPU.



# Chapter 5.

## ACKNOWLEDGMENTS

Some of the cuDNN library routines were derived from code developed by others and are subject to the following:

### 5.1. University of Tennessee

Copyright (c) 2010 The University of Tennessee.

All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- \* Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- \* Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer listed in this license in the documentation and/or other materials provided with the distribution.
- \* Neither the name of the copyright holders nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

## 5.2. University of California, Berkeley

### COPYRIGHT

All contributions by the University of California:  
Copyright (c) 2014, The Regents of the University of California (Regents)  
All rights reserved.

All other contributions:  
Copyright (c) 2014, the respective contributors  
All rights reserved.

Caffe uses a shared copyright model: each contributor holds copyright over their contributions to Caffe. The project versioning records all such contribution and copyright details. If a contributor wants to further mark their specific copyright on a particular contribution, they should indicate their copyright solely in the commit message of the change when it is committed.

### LICENSE

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

1. Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
2. Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

### CONTRIBUTION AGREEMENT

By contributing to the BVLC/caffe repository through pull-request, comment, or otherwise, the contributor releases their content to the license and copyright terms herein.

## **Notice**

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.

## **Trademarks**

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.

## **Copyright**

© 2007-2014 NVIDIA Corporation. All rights reserved.