

**NAME**

Stream Management –

**Typedefs**

```
typedef void(CUDART_CB * cudaStreamCallback_t)(cudaStream_t stream, cudaError_t status,
void *userData)
```

**Functions**

```
cudaError_t cudaStreamAddCallback (cudaStream_t stream, cudaStreamCallback_t callback,
void *userData, unsigned int flags)
```

*Add a callback to a compute stream.*

```
__cuda_builtin__ cudaError_t cudaStreamAttachMemAsync (cudaStream_t stream, void
*devPtr, size_t length, unsigned int flags)
```

*Attach memory to a stream asynchronously.*

```
cudaError_t cudaStreamCreate (cudaStream_t *pStream)
```

*Create an asynchronous stream.*

```
__cuda_builtin__ cudaError_t cudaStreamCreateWithFlags (cudaStream_t *pStream, unsigned
int flags)
```

*Create an asynchronous stream.*

```
__cuda_builtin__ cudaError_t cudaStreamCreateWithPriority (cudaStream_t *pStream,
unsigned int flags, int priority)
```

*Create an asynchronous stream with the specified priority.*

```
__cuda_builtin__ cudaError_t cudaStreamDestroy (cudaStream_t stream)
```

*Destroys and cleans up an asynchronous stream.*

```
__cuda_builtin__ cudaError_t cudaStreamGetFlags (cudaStream_t hStream, unsigned int *flags)
```

*Query the flags of a stream.*

```
__cuda_builtin__ cudaError_t cudaStreamGetPriority (cudaStream_t hStream, int *priority)
```

*Query the priority of a stream.*

```
cudaError_t cudaStreamQuery (cudaStream_t stream)
```

*Queries an asynchronous stream for completion status.*

```
cudaError_t cudaStreamSynchronize (cudaStream_t stream)
```

*Waits for stream tasks to complete.*

```
__cuda_builtin__ cudaError_t cudaStreamWaitEvent (cudaStream_t stream, cudaEvent_t event,
unsigned int flags)
```

*Make a compute stream wait on an event.***Detailed Description***\brief stream management functions of the CUDA runtime API (cuda\_runtime\_api.h)*

This section describes the stream management functions of the CUDA runtime application programming interface.

**Typedef Documentation**

```
typedef void(CUDART_CB * cudaStreamCallback_t)(cudaStream_t stream, cudaError_t status, void
*userData)
```

Type of stream callback functions.

**Parameters:***stream* The stream as passed to **cudaStreamAddCallback**, may be NULL.*status* **cudaSuccess** or any persistent error on the stream.*userData* User parameter provided at registration.**Function Documentation**

```
cudaError_t cudaStreamAddCallback (cudaStream_t stream, cudaStreamCallback_t callback, void *
userData, unsigned int flags)
```

Adds a callback to be called on the host after all currently enqueued items in the stream have completed. For each **cudaStreamAddCallback** call, a callback will be executed exactly once. The callback will block later work in the stream until it is finished.

The callback may be passed **cudaSuccess** or an error code. In the event of a device error, all subsequently executed callbacks will receive an appropriate **cudaError\_t**.

Callbacks must not make any CUDA API calls. Attempting to use CUDA APIs will result in **cudaErrorNotPermitted**. Callbacks must not perform any synchronization that may depend on



outstanding device work or other callbacks that are not mandated to run earlier. Callbacks without a mandated order (in independent streams) execute in undefined order and may be serialized.

This API requires compute capability 1.1 or greater. See **cudaDeviceGetAttribute** or **cudaGetDeviceProperties** to query compute capability. Calling this API with an earlier compute version will return **cudaErrorNotSupported**.

For the purposes of Unified Memory, callback execution makes a number of guarantees:

- The callback stream is considered idle for the duration of the callback. Thus, for example, a callback may always use memory attached to the callback stream.
- The start of execution of a callback has the same effect as synchronizing an event recorded in the same stream immediately prior to the callback. It thus synchronizes streams which have been 'joined' prior to the callback.
- Adding device work to any stream does not have the effect of making the stream active until all preceding callbacks have executed. Thus, for example, a callback might use global attached memory even if work has been added to another stream, if it has been properly ordered with an event.
- Completion of a callback does not cause a stream to become active except as described above. The callback stream will remain idle if no device work follows the callback, and will remain idle across consecutive callbacks without device work in between. Thus, for example, stream synchronization can be done by signaling from a callback at the end of the stream.

#### Parameters:

*stream* - Stream to add callback to

*callback* - The function to call once preceding stream operations are complete

*userData* - User specified data to be passed to the callback function

*flags* - Reserved for future use, must be 0

#### Returns:

**cudaSuccess**, **cudaErrorInvalidResourceHandle**, **cudaErrorNotSupported**

#### Note:

This function uses standard semantics.

Note that this function may also return error codes from previous, asynchronous launches.

#### See also:

**cudaStreamCreate**, **cudaStreamCreateWithFlags**, **cudaStreamQuery**,  
**cudaStreamSynchronize**, **cudaStreamWaitEvent**, **cudaStreamDestroy**, **cudaMallocManaged**,  
**cudaStreamAttachMemAsync**

**\_\_cuda\_builtin\_\_ cudaError\_t cudaStreamAttachMemAsync (cudaStream\_t stream, void \* devPtr, size\_t length, unsigned int flags)**

Enqueues an operation in *stream* to specify stream association of *length* bytes of memory starting from *devPtr*. This function is a stream-ordered operation, meaning that it is dependent on, and will only take effect when, previous work in *stream* has completed. Any previous association is automatically replaced.

*devPtr* must point to an address within managed memory space declared using the **\_\_managed\_\_** keyword or allocated with **cudaMallocManaged**.

*length* must be zero, to indicate that the entire allocation's stream association is being changed. Currently, it's not possible to change stream association for a portion of an allocation.

The stream association is specified using *flags* which must be one of **cudaMemAttachGlobal**, **cudaMemAttachHost** or **cudaMemAttachSingle**. If the **cudaMemAttachGlobal** flag is specified, the memory can be accessed by any stream on any device. If the **cudaMemAttachHost** flag is specified, the program makes a guarantee that it won't access the memory on the device from any stream. If the **cudaMemAttachSingle** flag is specified, the program makes a guarantee that it will only access the memory on the device from *stream*. It is illegal to attach singly to the NULL stream, because the NULL stream is a virtual global stream and not a specific stream. An error will be returned in this case. When memory is associated with a single stream, the Unified Memory system will allow CPU access to this memory region so long as all operations in *stream* have completed, regardless of whether other streams are active. In effect, this constrains exclusive ownership of the managed memory region by an active GPU to per-stream activity instead of whole-GPU activity.

Accessing memory on the device from streams that are not associated with it will produce undefined results. No error checking is performed by the Unified Memory system to ensure that kernels launched into other streams do not access this region.

It is a program's responsibility to order calls to **cudaStreamAttachMemAsync** via events,



synchronization or other means to ensure legal access to memory at all times. Data visibility and coherency will be changed appropriately for all kernels which follow a stream-association change. If `stream` is destroyed while data is associated with it, the association is removed and the association reverts to the default visibility of the allocation as specified at `cudaMallocManaged`. For `__managed__` variables, the default association is always `cudaMemAttachGlobal`. Note that destroying a stream is an asynchronous operation, and as a result, the change to default association won't happen until all work in the stream has completed.

**Parameters:**

*stream* - Stream in which to enqueue the attach operation  
*devPtr* - Pointer to memory (must be a pointer to managed memory)  
*length* - Length of memory (must be zero)  
*flags* - Must be one of `cudaMemAttachGlobal`, `cudaMemAttachHost` or `cudaMemAttachSingle`

**Returns:**

`cudaSuccess`, `cudaErrorNotReady`, `cudaErrorInvalidValue`  
`cudaErrorInvalidResourceHandle`

**Note:**

Note that this function may also return error codes from previous, asynchronous launches.

**See also:**

`cudaStreamCreate`, `cudaStreamCreateWithFlags`, `cudaStreamWaitEvent`,  
`cudaStreamSynchronize`, `cudaStreamAddCallback`, `cudaStreamDestroy`,  
`cudaMallocManaged`

**`cudaError_t cudaStreamCreate (cudaStream_t * pStream)`**

Creates a new asynchronous stream.

**Parameters:**

*pStream* - Pointer to new stream identifier

**Returns:**

`cudaSuccess`, `cudaErrorInvalidValue`

**Note:**

Note that this function may also return error codes from previous, asynchronous launches.

**See also:**

`cudaStreamCreateWithPriority`, `cudaStreamCreateWithFlags`, `cudaStreamGetPriority`,  
`cudaStreamGetFlags`, `cudaStreamQuery`, `cudaStreamSynchronize`, `cudaStreamWaitEvent`,  
`cudaStreamAddCallback`, `cudaStreamDestroy`

**`__cuda_builtin__ cudaError_t cudaStreamCreateWithFlags (cudaStream_t * pStream, unsigned int flags)`**

Creates a new asynchronous stream. The `flags` argument determines the behaviors of the stream.

Valid values for `flags` are

- **`cudaStreamDefault`**: Default stream creation flag.
- **`cudaStreamNonBlocking`**: Specifies that work running in the created stream may run concurrently with work in stream 0 (the NULL stream), and that the created stream should perform no implicit synchronization with stream 0.

**Parameters:**

*pStream* - Pointer to new stream identifier

*flags* - Parameters for stream creation

**Returns:**

`cudaSuccess`, `cudaErrorInvalidValue`

**Note:**

Note that this function may also return error codes from previous, asynchronous launches.

**See also:**

`cudaStreamCreate`, `cudaStreamCreateWithPriority`, `cudaStreamGetFlags`,  
`cudaStreamQuery`, `cudaStreamSynchronize`, `cudaStreamWaitEvent`,  
`cudaStreamAddCallback`, `cudaStreamDestroy`

**`__cuda_builtin__ cudaError_t cudaStreamCreateWithPriority (cudaStream_t * pStream, unsigned int flags, int priority)`**

Creates a stream with the specified priority and returns a handle in `pStream`. This API alters the scheduler priority of work in the stream. Work in a higher priority stream may preempt work already executing in a low priority stream.

`priority` follows a convention where lower numbers represent higher priorities. '0' represents



default priority. The range of meaningful numerical priorities can be queried using **cudaDeviceGetStreamPriorityRange**. If the specified priority is outside the numerical range returned by **cudaDeviceGetStreamPriorityRange**, it will automatically be clamped to the lowest or the highest number in the range.

**Parameters:**

*pStream* - Pointer to new stream identifier  
*flags* - Flags for stream creation. See **cudaStreamCreateWithFlags** for a list of valid flags that can be passed  
*priority* - Priority of the stream. Lower numbers represent higher priorities. See **cudaDeviceGetStreamPriorityRange** for more information about the meaningful stream priorities that can be passed.

**Returns:**

**cudaSuccess, cudaErrorInvalidValue**

**Note:**

Note that this function may also return error codes from previous, asynchronous launches. Stream priorities are supported only on Quadro and Tesla GPUs with compute capability 3.5 or higher. In the current implementation, only compute kernels launched in priority streams are affected by the stream's priority. Stream priorities have no effect on host-to-device and device-to-host memory operations.

**See also:**

**cudaStreamCreate, cudaStreamCreateWithFlags, cudaDeviceGetStreamPriorityRange, cudaStreamGetPriority, cudaStreamQuery, cudaStreamWaitEvent, cudaStreamAddCallback, cudaStreamSynchronize, cudaStreamDestroy**

**\_\_cuda\_builtin\_\_ cudaError\_t cudaStreamDestroy (cudaStream\_t stream)**

Destroys and cleans up the asynchronous stream specified by *stream*. In case the device is still doing work in the stream *stream* when **cudaStreamDestroy()** is called, the function will return immediately and the resources associated with *stream* will be released automatically once the device has completed all work in *stream*.

**Parameters:**

*stream* - Stream identifier

**Returns:**

**cudaSuccess, cudaErrorInvalidResourceHandle**

**Note:**

Note that this function may also return error codes from previous, asynchronous launches.

**See also:**

**cudaStreamCreate, cudaStreamCreateWithFlags, cudaStreamQuery, cudaStreamWaitEvent, cudaStreamSynchronize, cudaStreamAddCallback**

**\_\_cuda\_builtin\_\_ cudaError\_t cudaStreamGetFlags (cudaStream\_t hStream, unsigned int \* flags)**

Query the flags of a stream. The flags are returned in *flags*. See **cudaStreamCreateWithFlags** for a list of valid flags.

**Parameters:**

*hStream* - Handle to the stream to be queried  
*flags* - Pointer to an unsigned integer in which the stream's flags are returned

**Returns:**

**cudaSuccess, cudaErrorInvalidValue, cudaErrorInvalidResourceHandle**

**Note:**

Note that this function may also return error codes from previous, asynchronous launches.

**See also:**

**cudaStreamCreateWithPriority, cudaStreamCreateWithFlags, cudaStreamGetPriority**

**\_\_cuda\_builtin\_\_ cudaError\_t cudaStreamGetPriority (cudaStream\_t hStream, int \* priority)**

Query the priority of a stream. The priority is returned in *priority*. Note that if the stream was created with a priority outside the meaningful numerical range returned by **cudaDeviceGetStreamPriorityRange**, this function returns the clamped priority. See **cudaStreamCreateWithPriority** for details about priority clamping.

**Parameters:**

*hStream* - Handle to the stream to be queried  
*priority* - Pointer to a signed integer in which the stream's priority is returned

**Returns:**



**cudaSuccess, cudaErrorInvalidValue, cudaErrorInvalidResourceHandle**

**Note:**

Note that this function may also return error codes from previous, asynchronous launches.

**See also:**

**cudaStreamCreateWithPriority, cudaDeviceGetStreamPriorityRange, cudaStreamGetFlags**

**cudaError\_t cudaStreamQuery (cudaStream\_t stream)**

Returns **cudaSuccess** if all operations in *stream* have completed, or **cudaErrorNotReady** if not.

For the purposes of Unified Memory, a return value of **cudaSuccess** is equivalent to having called **cudaStreamSynchronize()**.

**Parameters:**

*stream* - Stream identifier

**Returns:**

**cudaSuccess, cudaErrorNotReady, cudaErrorInvalidResourceHandle**

**Note:**

Note that this function may also return error codes from previous, asynchronous launches.

**See also:**

**cudaStreamCreate, cudaStreamCreateWithFlags, cudaStreamWaitEvent, cudaStreamSynchronize, cudaStreamAddCallback, cudaStreamDestroy**

**\_\_cuda\_\_ builtin\_\_ cudaError\_t cudaStreamSynchronize (cudaStream\_t stream)**

Blocks until *stream* has completed all operations. If the **cudaDeviceScheduleBlockingSync** flag was set for this device, the host thread will block until the stream is finished with all of its tasks.

**Parameters:**

*stream* - Stream identifier

**Returns:**

**cudaSuccess, cudaErrorInvalidResourceHandle**

**Note:**

Note that this function may also return error codes from previous, asynchronous launches.

**See also:**

**cudaStreamCreate, cudaStreamCreateWithFlags, cudaStreamQuery, cudaStreamWaitEvent, cudaStreamAddCallback, cudaStreamDestroy**

**\_\_cuda\_\_ builtin\_\_ cudaError\_t cudaStreamWaitEvent (cudaStream\_t stream, cudaEvent\_t event, unsigned int flags)**

Makes all future work submitted to *stream* wait until *event* reports completion before beginning execution. This synchronization will be performed efficiently on the device. The event *event* may be from a different context than *stream*, in which case this function will perform cross-device synchronization.

The stream *stream* will wait only for the completion of the most recent host call to **cudaEventRecord()** on *event*. Once this call has returned, any functions (including **cudaEventRecord()** and **cudaEventDestroy()**) may be called on *event* again, and the subsequent calls will not have any effect on *stream*.

If **cudaEventRecord()** has not been called on *event*, this call acts as if the record has already completed, and so is a functional no-op.

**Parameters:**

*stream* - Stream to wait

*event* - Event to wait on

*flags* - Parameters for the operation (must be 0)

**Returns:**

**cudaSuccess, cudaErrorInvalidResourceHandle**

**Note:**

This function uses standard semantics.

Note that this function may also return error codes from previous, asynchronous launches.

**See also:**

**cudaStreamCreate, cudaStreamCreateWithFlags, cudaStreamQuery, cudaStreamSynchronize, cudaStreamAddCallback, cudaStreamDestroy**

**Author**

Generated automatically by Doxygen from the source code.

