



## **CudaExtension**

# Table of contents

Components

---

Extension for CUDA operations.

- UUID: d63a98fa-7882-11eb-a917-b38f664f399c
- Version: 2.0.0
- Author: NVIDIA
- License: LICENSE

## Components

### **nvidia::gxf::CudaStream**

Holds and provides access to native `cudaStream_t`.

`nvidia::gxf::CudaStream` handle must be allocated by `nvidia::gxf::CudaStreamPool`. Its lifecycle is valid until explicitly recycled through `nvidia::gxf::CudaStreamPool.releaseStream()` or implicitly until `nvidia::gxf::CudaStreamPool` is deactivated.

You may call `stream()` to get the native `cudaStream_t` handle, and to submit GPU operations. After the submission, GPU takes over the input tensors/buffers and keeps them in use. To prevent host carelessly releasing these in-use buffers, CUDA Codelet needs to call `record(event, input_entity, sync_cb)` to extend `input_entity`'s lifecycle until GPU completely consumes it. Alternatively, you may call `record(event, event_destroy_cb)` for native `cudaEvent_t` operations and free in-use resource via `event_destroy_cb`.

It is required to have a `nvidia::gxf::CudaStreamSync` in the graph pipeline after all the CUDA operations. See more details in `nvidia::gxf::CudaStreamSync`

- Component ID: 5683d692-7884-11eb-9338-c3be62d576be
- Defined in: `gxf/cuda/cuda_stream.hpp`

### **nvidia::gxf::CudaStreamId**

Holds CUDA stream Id to deduce `nvidia::gxf::CudaStream` handle.

`stream_cid` should be `nvidia::gxf::CudaStream` component id.

- Component ID: 7982aeac-37f1-41be-ade8-6f00b4b5d47c
- Defined in: `gxf/cuda/cuda_stream_id.hpp`

## **nvidia::gxf::CudaEvent**

Holds and provides access to native `cudaEvent_t` handle.

When a `nvidia::gxf::CudaEvent` is created, you'll need to initialize a native `cudaEvent_t` through `init(flags, dev_id)`, or set third party event through `initWithEvent(event, dev_id, free_fnc)`. The event keeps valid until `deinit` is called explicitly otherwise gets recycled in destructor.

- Component ID: f5388d5c-a709-47e7-86c4-171779bc64f3
- Defined in: `gxf/cuda/cuda_event.hpp`

## **nvidia::gxf::CudaStreamPool**

`CudaStream` allocation.

You must explicitly call `allocateStream()` to get a valid `nvidia::gxf::CudaStream` handle. This component would hold all the its allocated `nvidia::gxf::CudaStream` entities until `releaseStream(stream)` is called explicitly or the `CudaStreamPool` component is deactivated.

- Component ID: 6733bf8b-ba5e-4fae-b596-af2d1269d0e7
- Base Type: `nvidia::gxf::Allocator`

### **Parameters**

#### **dev\_id**

GPU device id.

- Flags: `GXF_PARAMETER_FLAGS_NONE`

- Type: GXF\_PARAMETER\_TYPE\_INT32
- Default Value: 0

### **stream\_flags**

Flag values to create CUDA streams.

- Flags: GXF\_PARAMETER\_FLAGS\_NONE
- Type: GXF\_PARAMETER\_TYPE\_INT32
- Default Value: 0

### **stream\_priority**

Priority values to create CUDA streams.

- Flags: GXF\_PARAMETER\_FLAGS\_NONE
- Type: GXF\_PARAMETER\_TYPE\_INT32
- Default Value: 0

### **reserved\_size**

User-specified file name without extension.

- Flags: GXF\_PARAMETER\_FLAGS\_NONE
- Type: GXF\_PARAMETER\_TYPE\_INT32

- Default Value: 1

## **max\_size**

Maximum Stream Size.

- Flags: GXF\_PARAMETER\_FLAGS\_NONE
- Type: GXF\_PARAMETER\_TYPE\_INT32
- Default Value: 0, no limitation.

## **nvidia::gxf::CudaStreamSync**

Synchronize all CUDA streams which are carried by message entities.

This codelet is required to get connected in the graph pipeline after all CUDA ops codelets. When a message entity is received, it would find all of the

`nvidia::gxf::CudaStreamId` in that message, and extract out each `nvidia::gxf::CudaStream`. With each `CudaStream` handle, it synchronizes all previous `nvidia::gxf::CudaStream.record()` events, along with all submitted GPU operations before this point.

### **Note**

`CudaStreamSync` must be set in the graph when `nvidia::gxf::CudaStream.record()` is used, otherwise it may cause memory leak.

- Component ID: 0d1d8142-6648-485d-97d5-277eed00129c
- Base Type: `nvidia::gxf::Codelet`

## **Parameters**

**rx**

Receiver to receive all messages carrying `nvidia::gxf::CudaStreamId` .

- Flags: `GXF_PARAMETER_FLAGS_NONE`
- Type: `GXF_PARAMETER_TYPE_HANDLE`
- Handle Type: `nvidia::gxf::Receiver`

## **tx**

Transmitter to send messages to downstream.

- Flags: `GXF_PARAMETER_FLAGS_OPTIONAL`
- Type: `GXF_PARAMETER_TYPE_HANDLE`
- Handle Type: `nvidia::gxf::Transmitter`

© Copyright 2022-2024, NVIDIA.. PDF Generated on 06/06/2024