



CUDA FOR TEGRA

DA-06762-001_v10.0 | March 2019

Application Note



TABLE OF CONTENTS

Chapter 1. Overview.....	1
--------------------------	---

Chapter 1.

OVERVIEW

This document provides an overview of NVIDIA® Tegra® memory architecture and considerations for porting code from a discrete GPU (dGPU) attached to an x86 system to the Tegra® integrated GPU (iGPU). It also discusses EGL interoperability.

This guide is for developers who are already familiar with programming in CUDA, and C/C++, and who want to develop applications for the Tegra® SoC.

Performance guidelines, best practices, terminology, and general information provided in the *CUDA C Programming Guide* and the *CUDA C Best Practices Guide* are applicable to all CUDA-capable GPU architectures, including Tegra® devices.

The *CUDA C Programming Guide* and the *CUDA C Best Practices Guide* are available at the following web sites:

CUDA C Programming Guide:

<http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>

CUDA C Best Practices Guide:

<http://docs.nvidia.com/cuda/cuda-c-best-practices-guide/index.html>

Chapter 2. MEMORY MANAGEMENT

In Tegra® devices, both the CPU (Host) and the iGPU share SoC DRAM memory. A dGPU with separate DRAM memory can be connected to the Tegra® device over PCIe or NVLink.

An overview of a dGPU-connected Tegra® memory system is shown in [Figure 1](#).

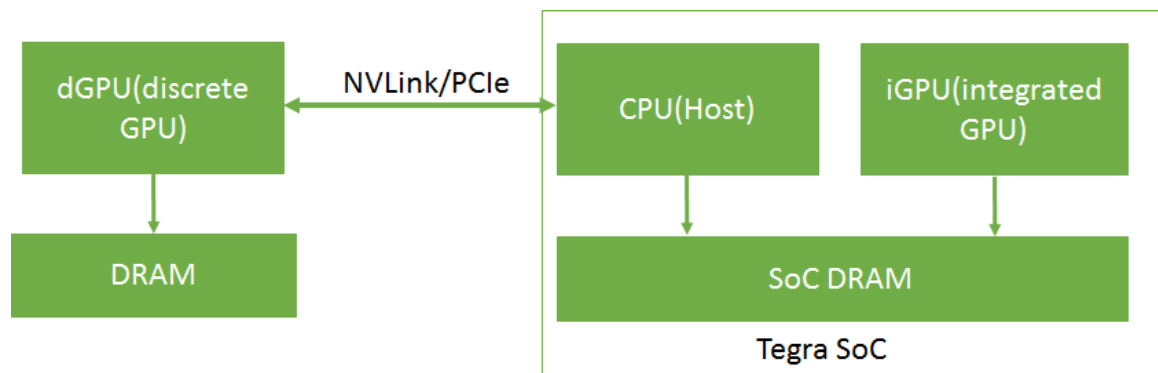


Figure 1 dGPU-connected Tegra Memory System

In Tegra®, device memory, host memory, and unified memory are allocated on the same physical SoC DRAM. On a dGPU, device memory is allocated on the dGPU DRAM. The caching behavior in a Tegra® system is different from that of an x86 system with a dGPU. The caching and accessing behavior of different memory types in a Tegra® system is shown in [Table 1](#).

Table 1 Characteristics of Different Memory Types in a Tegra System

Memory Type	CPU	iGPU	Tegra®-connected dGPU
Device memory	Not directly accessible	Cached	Cached
Pageable host memory	Cached	Not directly accessible	Not directly accessible
Pinned host memory	Uncached where compute capability is less than 7.2.	Uncached	Uncached

	Cached where compute capability is greater than or equal to 7.2.		
Unified memory	Cached	Cached	Not supported

On Tegra[®], because device memory, host memory, and unified memory are allocated on the same physical SoC DRAM, duplicate memory allocations and data transfers can be avoided.

Chapter 3.

PORTING CONSIDERATIONS

CUDA applications originally developed for dGPUs attached to x86 systems may require modifications to perform efficiently on Tegra® systems. This section describes the considerations for porting such applications to a Tegra® system, such as selecting an appropriate memory buffer type (pinned memory, unified memory, and others) and selecting between iGPU and dGPU, to achieve efficient performance for the application.

3.1. Memory Selection

CUDA applications can use various kinds of memory buffers, such as device memory, pageable host memory, pinned memory, and unified memory. Even though these memory buffer types are allocated on the same physical device, each has different accessing and caching behaviors, as shown in [Table 1](#). It is important to select the most appropriate memory buffer type for efficient application execution.

Device Memory

Use device memory for buffers whose accessibility is limited to the iGPU. For example, in an application with multiple kernels, there may be buffers that are used only by the intermediate kernels of the application as input or output. These buffers are accessed only by the iGPU. Such buffers should be allocated with device memory.

Pageable Host Memory

Use pageable host memory for buffers whose accessibility is limited to the CPU.

Pinned Memory

Tegra® systems with different compute capabilities exhibit different behavior in terms of I/O coherency. For example, Tegra® systems with compute capability greater than or equal to 7.2 are I/O coherent and others are not I/O coherent. On Tegra® systems with I/O coherency, the CPU access time of pinned memory is as good as pageable host memory because it is cached on the CPU. However, on Tegra® systems without I/O

coherency, the CPU access time of pinned memory is higher, because it is not cached on the CPU.

Pinned memory is recommended for small buffers because the caching effect is negligible for such buffers and also because pinned memory does not involve any additional overhead, unlike Unified Memory. With no additional overhead, pinned memory is also preferable for large buffers if the access pattern is not cache friendly on iGPU. For large buffers, when the buffer is accessed only once on iGPU in a coalescing manner, performance on iGPU can be as good as unified memory on iGPU.

Unified Memory

Unified memory is cached on the iGPU and the CPU. On Tegra®, using unified memory in applications requires additional coherency and cache maintenance operations during the kernel launch, synchronization and prefetching hint calls. This coherency maintenance overhead is slightly higher on a Tegra® system with compute capability less than 7.2 as they lack I/O coherency.

On Tegra® devices with I/O coherency (with a compute capability of 7.2 or greater) where unified memory is cached on both CPU and iGPU, for large buffers which are frequently accessed by the iGPU and the CPU and *the accesses on iGPU are repetitive*, unified memory is preferable since repetitive accesses can offset the cache maintenance cost. On Tegra® devices without I/O coherency (with a compute capability of less than 7.2), for large buffers which are frequently accessed by the CPU and the iGPU and *the accesses on iGPU are not repetitive*, unified memory is still preferable over pinned memory because pinned memory is not cached on both CPU and iGPU. That way, the application can take advantage of unified memory caching on the CPU.

Pinned memory or unified memory can be used to reduce the data transfer overhead between CPU and iGPU as both memories are directly accessible from the CPU and the iGPU. In an application, input and output buffers that must be accessible on both the host and the iGPU can be allocated using either unified memory or pinned memory.



The unified memory model requires the driver and system software to manage coherence on the current Tegra SOC. Software managed coherence is by nature non-deterministic and not recommended in a safe context. Zero-copy memory (pinned memory) is preferable in these applications.

Evaluate the impact of unified memory overheads, pinned memory cache misses, and device memory data transfers in applications to determine the correct memory selection.

3.2. Pinned Memory

This section provides guidelines for porting applications that use pinned memory allocations in x86 systems with dGPUs to Tegra®. CUDA applications developed for a

dGPU attached to x86 system use pinned memory to reduce data transfer time and to overlap data transfers with kernel execution time. For specific information on this topic, see “Data Transfer Between Host and Device” and “Asynchronous and Overlapping Transfers with Computation” at the following websites.

“Data Transfer Between Host and Device”:

<http://docs.nvidia.com/cuda/cuda-c-best-practices-guide/index.html#data-transfer-between-host-and-device>

“Asynchronous and Overlapping Transfers with Computation”:

<http://docs.nvidia.com/cuda/cuda-c-best-practices-guide/index.html#asynchronous-transfers-and-overlapping-transfers-with-computation>

On Tegra® systems with no I/O coherency, repetitive access of pinned memory degrades application performance, because pinned memory is not cached on the CPU in such systems.

A sample application is shown below in which a set of filters and operations (k1, k2, and k3) are applied to an image. Pinned memory is allocated to reduce data transfer time on an x86 system with a dGPU, increasing the overall application speed. However, targeting a Tegra® device with the same code causes a drastic increase in the execution time of the **readImage()** function because it repeatedly accesses an uncached buffer. This increases the overall application time. If the time taken by **readImage()** is significantly higher compared to kernels execution time, it is recommended to use unified memory to reduce the **readImage()** time. Otherwise, evaluate the application with pinned memory and unified memory by removing unnecessary data transfer calls to decide best suited memory.


```

// Sample code for an x86 system with a discrete GPU
int main()
{
    int *h_a,*d_a,*d_b,*d_c,*d_d,*h_d;
    int height = 1024;
    int width = 1024;
    size_t sizeOfImage = width * height * sizeof(int); // 4MB image

    //Pinned memory allocated to reduce data transfer time
    cudaMallocHost(h_a, sizeOfImage);
    cudaMallocHost(h_d, sizeOfImage);

    //Allocate buffers on GPU
    cudaMalloc(&d_a, sizeOfImage);
    cudaMalloc(&d_b, sizeOfImage);
    cudaMalloc(&d_c, sizeOfImage);
    cudaMalloc(&d_d, sizeOfImage);

    //CPU reads Image;
    readImage(h_a); // Intialize the h_a buffer

    // Transfer image to GPU
    cudaMemcpy(d_a, h_a, sizeOfImage, cudaMemcpyHostToDevice);

    // Data transfer is fast as we used pinned memory
    // ----- CUDA Application pipeline start -----
    k1<<<...>>>(d_a,d_b) // Apply filter 1
    k2<<<...>>>(d_b,d_c) // Apply filter 2
    k3<<<...>>>(d_c,d_d) // Some operation on image data
    // ----- CUDA Application pipeline end -----

    // Transfer processed image to CPU
    cudaMemcpy(h_d, d_d, sizeOfImage, cudaMemcpyDeviceToHost);
    // Data transfer is fast as we used pinned memory

    // Use processed Image i.e h_d in later computations on CPU.
    UseImageonCPU(h_d);
}

// Porting the code on Tegra
int main()
{
    int *h_a,*d_b,*d_c,*h_d;
    int height = 1024;
    int width = 1024;
    size_t sizeOfImage = width * height * sizeof(int); // 4MB image

    //Unified memory allocated for input and output
    //buffer of application pipeline
    cudaMallocManaged(h_a, sizeOfImage,cudaMemAttachHost);
    cudaMallocManaged(h_d, sizeOfImage);

    //Intermediate buffers not needed on CPU side.
    //So allocate them on device memory
    cudaMalloc(&d_b, sizeOfImage);
    cudaMalloc(&d_c, sizeOfImage);

    //CPU reads Image;
    readImage(h_a); // Intialize the h_a buffer
    // ----- CUDA Application pipeline start -----
    // Prefetch input image data to GPU
    cudaStreamAttachMemAsync(NULL, h_a, 0, cudaMemAttachGlobal);
    k1<<<...>>>(h_a,d_b)
    k2<<<...>>>(d_b,d_c)
    k3<<<...>>>(d_c,h_d)
    // Prefetch output image data to CPU
    cudaStreamAttachMemAsync(NULL, h_d, 0, cudaMemAttachHost);
    cudaStreamSynchronize(NULL);
    // ----- CUDA Application pipeline end -----

    // Use processed Image i.e h_d on CPU side.
    UseImageonCPU(h_d);
}

```

The `cudaHostRegister()` function

The `cudaHostRegister()` function is not supported on Tegra® devices with compute capability less than 7.2, because those devices do not have I/O coherency. Use other pinned memory allocation functions such as `cudaMallocHost()` and `cudaHostAlloc()` if `cudaHostRegister()` is not supported on the device.

GNU Atomic operations on pinned memory

The GNU atomic operations on uncached memory is not supported on Tegra® CPU. As pinned memory is not cached on Tegra® devices with compute capability less than 7.2, GNU atomic operations is not supported on pinned memory.

3.3. Effective Usage of Unified Memory on Tegra

Using unified memory in applications requires additional coherency and cache maintenance operations at kernel launch, synchronization, and prefetching hint calls. These operations are performed synchronously with other GPU work which can cause unpredictable latencies in the application.

The performance of unified memory on Tegra® can be improved by providing data prefetching hints. The driver can use these prefetching hints to optimize the coherence operations. To prefetch the data, the `cudaStreamAttachMemAsync()` function can be used, in addition to the techniques described in the “Coherency and Concurrency” section of the *CUDA C Programming Guide* at the following link:

<http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#um-coherency-hd>

to prefetch the data. The prefetching behavior of unified memory, as triggered by the changing states of the attachment flag, is shown in Table 2.

Table 2 Unified Memory Prefetching Behavior per Changing Attachment Flag States

Previous Flag	Current Flag	Prefetching Behavior
<code>cudaMemAttachGlobal/</code> <code>cudaMemAttachSingle</code>	<code>cudaMemAttachHost</code>	Causes prefetch to CPU
<code>cudaMemAttachHost</code>	<code>cudaMemAttachGlobal/</code> <code>cudaMemAttachSingle</code>	Causes prefetch to GPU
<code>cudaMemAttachGlobal</code>	<code>cudaMemAttachSingle</code>	No prefetch to GPU
<code>cudaMemAttachSingle</code>	<code>cudaMemAttachGlobal</code>	No prefetch to GPU

The following example shows usage of `cudaStreamAttachMemAsync()` to prefetch data.



However, not supported on Tegra® devices are the data prefetching techniques that use `cudaMemPrefetchAsync()` as described in the “Performance Tuning” section of the *CUDA C Programming Guide* at the following web site:

<http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#um-performance-tuning>



There are limitations in QNX system software which prevent implementation of all UVM optimizations. Because of this, using `cudaStreamAttachMemAsync()` to prefetch hints on QNX does not benefit performance.

```
__global__ void matrixMul(int *p, int *q, int *r, int hp, int hq, int wp, int wq)
{
    // Matrix multiplication kernel code
}
void MatrixMul(int hp, int hq, int wp, int wq)
{
    int *p, *q, *r;
    int i;
    size_t sizeP = hp*wp*sizeof(int);
    size_t sizeQ = hq*wq*sizeof(int);
    size_t sizeR = hp*wq*sizeof(int);

    //Attach buffers 'p' and 'q' to CPU and buffer 'r' to GPU
    cudaMallocManaged(&p, sizeP, cudaMemAttachHost);
    cudaMallocManaged(&q, sizeQ, cudaMemAttachHost);
    cudaMallocManaged(&r, sizeR);
    //Intialize with random values
    randFill(p, q, hp, wp, hq, wq);

    // Prefetch p,q to GPU as they are needed in computation
    cudaStreamAttachMemAsync(NULL, p, 0, cudaMemAttachGlobal);
    cudaStreamAttachMemAsync(NULL, q, 0, cudaMemAttachGlobal);
    matrixMul<<<...>>>(p, q, r, hp, hq, wp, wq);

    // Prefetch 'r' to CPU as only 'r' is needed
    cudaStreamAttachMemAsync(NULL, r, 0, cudaMemAttachHost);
    cudaStreamSynchronize(NULL);

    // Print buffer 'r' values
    for(i = 0; i < hp*wq; i++)
        printf("%d ", r[i]);
}
```



An additional `cudaStreamSynchronize(NULL)` call can be added after the `matrixMul` kernel code to avoid callback threads that cause unpredictability in a `cudaStreamAttachMemAsync()` call.

3.4. GPU Selection

On a Tegra system with a dGPU, deciding whether a CUDA application runs on the iGPU or the dGPU can have implications for the performance of the application. Some of the factors that need to be considered while making such a decision are kernel execution

time, data transfer time, data locality, and latency. For example, to run an application on a dGPU, data must be transferred between the SoC and the dGPU. This data transfer can be avoided if the application runs on an iGPU.

3.5. Synchronization Mechanism Selection

The `cudaSetDeviceFlags` API is used to control the synchronization behaviour of CPU thread. Until CUDA 10.1, by default, the synchronization mechanism on iGPU uses `cudaDeviceBlockingSync` flag, which blocks the CPU thread on a synchronization primitive when waiting for the device to finish work. The `cudaDeviceBlockingSync` flag is suited for platforms with power constraints. But on platforms which requires low latency, `cudaDeviceScheduleSpin` flag needs to be set manually. Since CUDA 10.1, for each platform, the default synchronization flag is determined based on what is optimized for that platform. More information about the synchronization flags is given at `cudaSetDeviceFlags` API documentation.

3.6. CUDA Features Not Supported on Tegra

All core features of CUDA are supported on Tegra platforms. The exceptions are listed below.

- ▶ The `cudaHostRegister()` function is not supported on QNX systems. This is due to the limitations on QNX OS. It is supported in Linux systems with compute capability greater than or equal to 7.2.
- ▶ System wide atomics are not supported on Tegra devices with compute capability less than 7.2.
- ▶ Unified memory is not supported on dGPU attached to Tegra.
- ▶ `cudaMemPrefetchAsync()` function is not supported since unified memory with concurrent access is not yet supported on iGPU.
- ▶ NVIDIA management library (NVML) library is not supported on Tegra. However, as an alternative to monitor the resource utilization, `tegrastats` can be used.
- ▶ CUDA IPC (CUDA Inter-process communication) is not supported on Tegra devices. EGLStream can be used to communicate between CUDA contexts in two processes.
- ▶ Remote direct memory access (RDMA) is not supported on Tegra devices.
- ▶ JIT compilation might require a considerable amount of CPU and bandwidth resources, potentially interfering with other workloads in the system. Thus, JIT compilations such as PTX-JIT and NVRTC JIT are not recommended for deterministic automotive applications and can be bypassed completely by compiling for specific GPU targets. JIT compilation is not supported on Tegra devices in the safe context.
- ▶ Multi process service (MPS) is not supported on Tegra.
- ▶ Peer to peer (P2P) communication calls are not supported on Tegra.
- ▶ The cuSOLVER library is not supported on Tegra® systems running QNX.
- ▶ The nvGRAPH library is not supported.

More information on some of these features can be found at the following web sites:

IPC:

<http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#interprocess-communication>

RDMA:

<http://docs.nvidia.com/cuda/gpudirect-rdma/index.html>

MPS:

https://docs.nvidia.com/deploy/pdf/CUDA_Multi_Process_Service_Overview.pdf

P2P:

<http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#peer-to-peer-memory-access>

Chapter 4.

EGL INTEROPERABILITY

An interop is an efficient mechanism to share resources between two APIs. To share data with multiple APIs, an API must implement an individual interop for each.

EGL provides interop extensions that allow it to function as a hub connecting APIs, removing the need for multiple interops, and encapsulating the shared resource. An API must implement these extensions to interoperate with any other API via EGL. The CUDA supported EGL interops are EGLStream, EGLImage, and EGLSync.

EGL interop extensions allow applications to switch between APIs without the need to rewrite code. For example, an EGLStream-based application in which *NvMedia* is the producer and *CUDA* is the consumer can be modified to use *OpenGL* as the consumer without modifying the producer code.

4.1. EGLStream

EGLStream interoperability facilitates efficient transfer of a sequence of frames from one API to another API, allowing use of multiple *Tegra*® engines such as CPU, GPU, ISP, and others.

Consider an application where a camera captures images continuously, shares them with *CUDA* for processing, and then later renders those images using *OpenGL*. In this application, the image frames are shared across *NvMedia*, *CUDA* and *OpenGL*. The absence of EGLStream interoperability would require the application to include multiple interops and redundant data transfers between APIs. EGLStream has one producer and one consumer.

EGLStream offers the following benefits:

- ▶ Efficient transfer of frames between a producer and a consumer.
- ▶ Implicit synchronization handling.
- ▶ Cross-process support.
- ▶ dGPU and iGPU support.
- ▶ Linux, QNX, and Android operating system support.

4.1.1. EGLStream Flow

The EGLStream flow has the following steps:

1. Initialize producer and consumer APIs
2. Create an EGLStream and connect the consumer and the producer.



EGLStream is created using `eglCreateStreamKHR()` and destroyed using `eglDestroyStreamKHR()`.

The consumer should always connect to EGLStream before the producer.

For more information see the EGLStream specification at the following web site:
https://www.khronos.org/registry/EGL/extensions/KHR/EGL_KHR_stream.txt

3. Allocate memory used for EGL frames.
4. The producer populates an EGL frame and presents it to EGLStream.
5. The consumer acquires the frame from EGLStream and releases it back to EGLStream after processing.
6. The producer collects the consumer-released frame from EGLStream.
7. The producer presents the same frame, or a new frame to EGLStream.
8. Steps 4-7 are repeated until completion of the task, with an old frame or a new frame.
9. The consumer and the producer disconnect from EGLStream.
10. Deallocate the memory used for EGL frames.
11. De-initialize the producer and consumer APIs.

EGLStream application flow is shown in Figure 2.

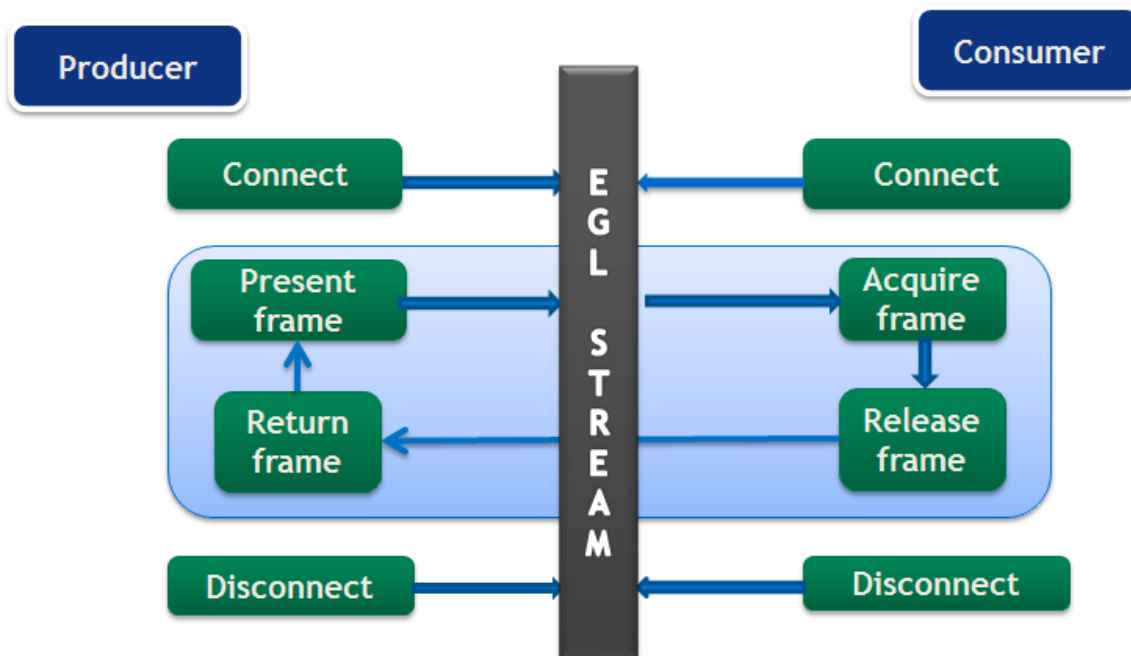


Figure 2 EGLStream Flow

CUDA producer and consumer functions are listed in [Table 3](#).

Table 3 CUDA Producer and Consumer Functions

Role	Functionality	API
Producer	To connect a producer to EGLStream	cuEGLStreamProducerConnect() cudaEGLStreamProducerConnect()
	To present frame to EGLStream	cuEGLStreamProducerPresentFrame() cudaEGLStreamProducerPresentFrame()
	Obtain released frames	cuEGLStreamProducerReturnFrame() cudaEGLStreamProducerReturnFrame()
	To disconnect from EGLStream	cuEGLStreamProducerDisconnect() cudaEGLStreamProducerDisconnect()
Consumer	To connect a consumer to EGLStream	cuEGLStreamConsumerConnect() cuEGLStreamConsumeConnectWithFlags() cudaEGLStreamConsumerConnect() cudaEGLStreamConsumerConnectWithFlags()
	To acquire frame from EGLStream	cuEGLStreamConsumerAcquireFrame() cudaEGLStreamConsumerAcquireFrame()
	To release the consumed frame	cuEGLStreamConsumerReleaseFrame() cudaEGLStreamConsumerReleaseFrame()
	To disconnect from EGLStream	cuEGLStreamConsumerDisconnect() cudaEGLStreamConsumerDisconnect()

4.1.2. CUDA as Producer

When CUDA is the producer, the supported consumers are CUDA, NvMedia and OpenGL. API functions to be used when CUDA is the producer are listed in [Table 3](#). Except for connecting and disconnecting from EGLStream, all API calls are non-blocking.

The following producer side steps are shown in the example code that follows:

1. Prepare a frame (lines 3-19).
2. Connect the producer to EGLStream (line 21).
3. Populate the frame and present to EGLStream (lines 23-25).
4. Get the released frame back from EGLStream (Line 27).

5. Disconnect the consumer after completion of the task. (Line 31).

```
void ProducerThread(EGLStreamKHR eglStream) {
    //Prepares frame
    cudaEglFrame* cudaEgl = (cudaEglFrame *)malloc(sizeof(cudaEglFrame));
    cudaEgl->planeDesc[0].width = WIDTH;
    cudaEgl->planeDesc[0].depth = 0;
    cudaEgl->planeDesc[0].height = HEIGHT;
    cudaEgl->planeDesc[0].numChannels = 4;
    cudaEgl->planeDesc[0].pitch = WIDTH * cudaEgl->planeDesc[0].numChannels;
    cudaEgl->frameType = cudaEglFrameTypePitch;
    cudaEgl->planeCount = 1;
    cudaEgl->eglColorFormat = cudaEglColorFormatARGB;
    cudaEgl->planeDesc[0].channelDesc.f=cudaChannelFormatKindUnsigned
    cudaEgl->planeDesc[0].channelDesc.w = 8;
    cudaEgl->planeDesc[0].channelDesc.x = 8;
    cudaEgl->planeDesc[0].channelDesc.y = 8;
    cudaEgl->planeDesc[0].channelDesc.z = 8;
    size_t numElem = cudaEgl->planeDesc[0].pitch * cudaEgl->planeDesc[0].height;
    // Buffer allocated by producer
    cudaMalloc(&(cudaEgl->pPitch[0].ptr), numElem);
    //CUDA producer connects to EGLStream
    cudaEGLStreamProducerConnect(&conn, eglStream, WIDTH, HEIGHT))
    // Sets all elements in the buffer to 1
    K1<<<...>>>(cudaEgl->pPitch[0].ptr, 1, numElem);
    // Present frame to EGLStream
    cudaEGLStreamProducerPresentFrame(&conn, *cudaEgl, NULL);

    cudaEGLStreamProducerReturnFrame(&conn, cudaEgl, eglStream);
    .
    .
    //clean up
    cudaEGLStreamProducerDisconnect(&conn);
    .
}
```

A frame is represented as a **cudaEglFrame** structure. The **frameType** parameter in **cudaEglFrame** indicates the memory layout of the frame. The supported memory layouts are CUDA Array and device pointer. Any mismatch in the width and height values of frame with the values specified in **cudaEGLStreamProducerConnect()** leads to undefined behavior. In the sample, the CUDA producer is sending a single frame, but it can send multiple frames over a loop. CUDA cannot present more than 64 active frames to EGLStream.

The **cudaEGLStreamProducerReturnFrame()** call waits until it receives the released frame from the consumer. Once the CUDA producer presents the first frame to EGLStream, at least one frame is always available for consumer acquisition until the producer disconnects. This prevents the removal of the last frame from EGLStream, which would block **cudaEGLStreamProducerReturnFrame()**.

Use the **EGL_NV_stream_reset** extension to set EGLStream attribute **EGL_SUPPORT_REUSE_NV** to false to allow the last frame to be removed from EGLStream. This allows removing or returning the last frame from EGLStream.

4.1.3. CUDA as Consumer

When CUDA is the consumer, the supported producers are CUDA, OpenGL, NvMedia, Argus, and Camera. API functions to be used when CUDA is the consumer are listed in

Table 3. Except for connecting and disconnecting from EGLStream, all API calls are non-blocking.

The following consumer side steps are shown in the sample code that follows:

1. Connect consumer to EGLStream (line 5).
2. Acquire frame from EGLStream (lines 8-10).
3. Process the frame on consumer (line 16).
4. Release frame back to EGLStream (line 19).
5. Disconnect the consumer after completion of the task (line 22).

```
void ConsumerThread(EGLStreamKHR eglStream) {
    .
    .
    //Connect consumer to EGLStream
    cudaEGLStreamConsumerConnect(&conn, eglStream);
    // consumer acquires a frame
    unsigned int timeout = 16000;
    cudaEGLStreamConsumerAcquireFrame(& conn, &cudaResource, eglStream, timeout);
    //consumer gets a cuda object pointer
    cudaGraphicsResourceGetMappedEglFrame(&cudaEgl, cudaResource, 0, 0);
    size_t numElem = cudaEgl->planeDesc[0].pitch * cudaEgl->planeDesc[0].height;
    .
    .
    int checkIfOne = 1;
    // Checks if each value in the buffer is 1, if any value is not 1, it sets
    // checkIfOne = 0.
    K2<<<...>>>(cudaEgl->pPitch[0].ptr, 1, numElem, checkIfOne);
    .
    .
    cudaEGLStreamConsumerReleaseFrame(&conn, cudaResource, &eglStream);
    .
    .
    cudaEGLStreamConsumerDisconnect(&conn);
    .
    .
}
```

In the sample code, the CUDA consumer receives a single frame, but it can also receive multiple frames over a loop. If a CUDA consumer fails to receive a new frame in the specified time limit using `cudaEGLStreamConsumerAcquireFrame()`, it reacquires the previous frame from EGLStream. The time limit is indicated by the timeout parameter.

The application can use `eglQueryStreamKHR()` to query for the availability of new frames using. If the consumer uses already released frames, it results in undefined behavior. The consumer behavior is defined only for read operations. Behavior is undefined when the consumer writes to a frame.

If the CUDA context is destroyed while connected to EGLStream, the stream is placed in the `EGL_STREAM_STATE_DISCONNECTED_KHR` state and the connection handle is invalidated.

4.1.4. Implicit Synchronization

EGLStream provides implicit synchronization in an application. For example, in the previous code samples, both the producer and consumer threads are running in parallel and the K1 and K2 kernel processes access the same frame, but K2 execution in the consumer thread is guaranteed to occur only after kernel K1 in the producer thread finishes. The `cudaEGLStreamConsumerAcquireFrame()` function waits on the GPU

side until K1 finishes and ensures synchronization between producer and consumer. The variable `checkIfOne` is never set to 0 inside the K2 kernel in the consumer thread.

Similarly, `cudaEGLStreamProducerReturnFrame()` in the producer thread is guaranteed to get the frame only after K2 finishes and the consumer releases the frame. These non-blocking calls allow the CPU to do other computation in between, as synchronization is taken care of on the GPU side.

The `EGLStreams_CUDA_Interop` CUDA sample code shows the usage of EGLStream in detail.

4.1.5. Data Transfer Between Producer and Consumer

Data transfer between producer and consumer is avoided when they are present on the same device. In a Tegra® platform that includes a dGPU however, such as is in NVIDIA DRIVE™ PX 2, the producer and consumer can be present on different devices. In that case, an additional memory copy is required internally to move the frame between Tegra® SoC DRAM and dGPU DRAM. EGLStream allows producer and consumer to run on any GPU without code modification.



On systems where a Tegra® device is connected to a dGPU, if a producer frame uses CUDA array, both producer and consumer should be on the same GPU. But if a producer frame uses CUDA device pointers, the consumer can be present on any GPU.

4.1.6. EGLStream Pipeline

An application can use multiple EGL streams in a pipeline to pass the frames from one API to another. For an application where NvMedia sends a frame to CUDA for computation, CUDA sends the same frame to OpenGL for rendering after the computation.

The EGLStream pipeline is illustrated in Figure 3.



Figure 3 EGLStream Pipeline

NvMedia and CUDA connect as producer and consumer respectively to one EGLStream. CUDA and OpenGL connect as producer and consumer respectively to another EGLStream.

Using multiple EGLStreams in pipeline fashion gives the flexibility to send frames across multiple APIs without allocating additional memory or requiring explicit data transfers. Sending a frame across the above EGLStream pipeline involves the following steps.

1. NvMedia sends a frame to CUDA for processing.
2. CUDA uses the frame for computation and sends to OpenGL for rendering.
3. OpenGL consumes the frame and releases it back to CUDA.
4. CUDA releases the frame back to NvMedia.

The above steps can be performed in a loop to facilitate the transfer of multiple frames in the EGLStream pipeline.

4.2. EGLImage

An EGLImage interop allows an EGL client API to share image data with other EGL client APIs. For example, an application can use an EGLImage interop to share an OpenGL texture with CUDA without allocating any additional memory. A single EGLImage object can be shared across multiple client APIs for modification.

An EGLImage interop does not provide implicit synchronization. Applications must maintain synchronization to avoid race conditions.



An EGLImage is created using `eglCreateImageKHR()` and destroyed using `eglDestroyImageKHR()`.

For more information see the EGLImage specification at the following web site:

https://www.khronos.org/registry/EGL/extensions/KHR/EGL_KHR_image_base.txt

4.2.1. CUDA interop with EGLImage

CUDA supports interoperation with EGLImage, allowing CUDA to read or modify the data of an EGLImage. An EGLImage can be a single or multi-planar resource. In CUDA, a single-planar EGLImage object is represented as a CUDA array or device pointer. Similarly, a multi-planar EGLImage object is represented as an array of device pointers or CUDA arrays. EGLImage is supported on Tegra® devices running the Linux, QNX, or Android operating systems.

Use the `cudaGraphicsEGLRegisterImage()` API to register an EGLImage object with CUDA. Registering an EGLImage with CUDA creates a graphics resource object. An application can use `cudaGraphicsResourceGetMappedEglFrame()` to get a frame from the graphics resource object. In CUDA, a frame is represented as a `cudaEglFrame` structure. The `frameType` parameter in `cudaEglFrame` indicates if the frame is a CUDA device pointer or a CUDA array. For a single planar graphics resource, an application can directly obtain a device pointer or CUDA array using `cudaGraphicsResourceGetMappedPointer()` or `cudaGraphicsSubResourceGetMappedArray()` respectively. A CUDA array can be bound to a texture or surface reference to access inside a kernel. Also, a multi-dimensional CUDA array can be read and written via `cudaMemcpy3D()`.



An EGLImage cannot be created from a CUDA object. The `cudaGraphicsEGLRegisterImage()` function is only supported on Tegra® devices.

The following sample code shows EGLImage interoperability. In the code, an EGLImage object `eglImage` is created using OpenGL texture. The `eglImage` object is mapped as a CUDA array `pArray` in CUDA. The `pArray` array is bound to a surface object to allow

modification of the OpenGL texture in the changeTexture. The function `checkBuf()` checks if the texture is updated with new values.

```
int width = 256;
int height = 256;
int main()
{
    .
    .
    unsigned char *hostSurf;
    unsigned char *pSurf;
    CUarray pArray;
    unsigned int bufferSize = WIDTH * HEIGHT * 4;
    pSurf= (unsigned char *)malloc(bufferSize); hostSurf = (unsigned char
    *)malloc(bufferSize);
    // Initialize the buffer
    for(int y = 0; y < HEIGHT; y++)
    {
        for(int x = 0; x < WIDTH; x++)
        {
            pSurf[(y*WIDTH + x) * 4] = 0; pSurf[(y*WIDTH + x) * 4 + 1] = 0;
            pSurf[(y*WIDTH + x) * 4 + 2] = 0; pSurf[(y*WIDTH + x) * 4 + 3] = 0;
        }
    }

    // NOP call to error-check the above glut calls
    GL_SAFE_CALL({});

    //Init texture
    GL_SAFE_CALL(glGenTextures(1, &tex));
    GL_SAFE_CALL(glBindTexture(GL_TEXTURE_2D, tex));
    GL_SAFE_CALL(glTexImage2D(GL_TEXTURE_2D, 0, GL_RGBA, WIDTH, HEIGHT, 0, GL_RGBA,
    GL_UNSIGNED_BYTE, pSurf));

    EGLDisplay eglDisplayHandle = eglGetCurrentDisplay();
    EGLContext eglCtx = eglGetCurrentContext();

    // Create the EGL_Image
    EGLint eglImgAttrs[] = { EGL_IMAGE_PRESERVED_KHR, EGL_FALSE, EGL_NONE,
    EGL_NONE };
    EGLImageKHR eglImage = eglCreateImageKHR(eglDisplayHandle, eglCtx,
    EGL_GL_TEXTURE_2D_KHR, (EGLClientBuffer)(intptr_t)tex, eglImgAttrs);
    glFinish();
    glTexSubImage2D(GL_TEXTURE_2D, 0, 0, 0, WIDTH, HEIGHT, GL_RGBA,
    GL_UNSIGNED_BYTE, pSurf);
    glFinish();

    // Register buffer with CUDA
    cuGraphicsEGLRegisterImage(&pResource, eglImage,
    CU_GRAPHICS_REGISTER_FLAGS_NONE);
    //Get CUDA array from graphics resource object
    cuGraphicsSubResourceGetMappedArray( &pArray, pResource, 0, 0);

    cuCtxSynchronize();

    //Create a CUDA surface object from pArray
    CUsresult status = CUDA_SUCCESS;
    CUDA_RESOURCE_DESC wdsc;
    memset(&wdsc, 0, sizeof(wdsc));
    wdsc.resType = CU_RESOURCE_TYPE_ARRAY; wdsc.res.array.hArray = pArray;
    CUsurfObject writeSurface;
    cuSurfObjectCreate(&writeSurface, &wdsc);

    dim3 blockSize(32,32);
    dim3 gridSize(width/blockSize.x,height/blockSize.y);
    // Modifies the OpenGL texture using CUDA surface object
    changeTexture<<<gridSize, blockSize>>>>(writeSurface, width, height);
    cuCtxSynchronize();

    CUDA_MEMCPY3D cpdesc;
    memset(&cpdesc, 0, sizeof(cpdesc));
    cpdesc.srcXInBytes = cpdesc.srcY = cpdesc.srcZ = cpdesc.srcLOD = 0;
    cpdesc.dstXInBytes = cpdesc.dstY = cpdesc.dstZ = cpdesc.dstLOD = 0;
    cpdesc.srcMemoryType = CU_MEMORYTYPE_ARRAY; cpdesc.dstMemoryType =
    CU_MEMORYTYPE_HOST;
    cpdesc.srcArray = pArray; cpdesc.dstHost = (void *)hostSurf;
```

Because EGLImage does not provide implicit synchronization, the above sample application uses `glFinish()` and `cudaThreadSynchronize()` calls to achieve synchronization. Both calls block the CPU thread. To avoid blocking the CPU thread, use EGLSync to provide synchronization. An example using EGLImage and EGLSync is shown in the following section.

4.3. EGLSync

EGLSync is a cross-API synchronization primitive. It allows an EGL client API to share its synchronization object with other EGL client APIs. For example, applications can use an EGLSync interop to share the OpenGL synchronization object with CUDA.



An EGLSync object is created using `eglCreateSyncKHR()` and destroyed using `eglDestroySyncKHR()`.

For more information see the EGLSync specification at the following web site:

https://www.khronos.org/registry/EGL/extensions/KHR/EGL_KHR_fence_sync.txt

4.3.1. CUDA Interop with EGLSync

In an imaging application, where two clients run on a GPU and share a resource, the absence of a cross-API GPU synchronization object forces the clients to use CPU-side synchronization to avoid race conditions. The CUDA interop with EGLSync allows the application to exchange synchronization objects between CUDA and other client APIs directly. This avoids the need for CPU-side synchronization and allows CPU to complete other tasks. In CUDA, an EGLSync object is mapped as a CUDA event.



Currently CUDA interop with EGLSync is supported only on Tegra® devices.

4.3.2. Creating EGLSync from a CUDA Event

Creating an EGLSync object from a CUDA event is shown in the following sample code.

```
EGLDisplay dpy = eglGetCurrentDisplay();
// Create CUDA event
cudaEvent_t event;
cudaStream_t *stream;
cudaEventCreate(&event);
cudaStreamCreate(&stream);
// Record the event with cuda event
cudaEventRecord(event, stream);
const EGLAttrib attribs[] = {
    EGL_CUDA_EVENT_HANDLE_NV, (EGLAttrib) event,
    EGL_NONE
};
//Create EGLSync from the cuda event
eglsync = eglCreateSync(dpy, EGL_NV_CUDA_EVENT_NV, attribs);
//Wait on the sync
eglWaitSyncKHR(...);
```



Initialize a CUDA event before creating an EGLSync object from it to avoid undefined behavior.

4.3.3. Creating a CUDA Event from EGLSync

Creating a CUDA event from an EGLSync object is shown in the following sample code.

```
EGLSync eglsync;
EGLDisplay dpy = eglGetCurrentDisplay();
// Create an eglSync object from OpenGL fence sync object
eglsync = eglCreateSyncKHR(dpy, EGL_SYNC_FENCE_KHR, NULL);
cudaEvent_t event;
cudaStream_t* stream;
cudaStreamCreate(&stream);
// Create CUDA event from eglSync
cudaEventCreateFromEGLSync(&event, eglsync, cudaEventDefault);
// Wait on the cuda event. It waits on GPU till OpenGL finishes its
// task
cudaStreamWaitEvent(stream, event, 0);
```



The `cudaEventRecord()` and `cudaEventElapsedTime()` functions are not supported for events created from an EGLSync object.

The same example given in the EGLImage section is re-written below to illustrate the usage of an EGLSync interop. In the sample code, the CPU blocking calls such as **glFinish()** and **cudaThreadSynchronize()** are replaced with EGLSync interop calls.

```
int width = 256;
int height = 256;
int main()
{
    .
    .
    unsigned char *hostSurf;
    unsigned char *pSurf;
    cudaArray_t pArray;
    unsigned int bufferSize = WIDTH * HEIGHT * 4;
    pSurf= (unsigned char *)malloc(bufferSize); hostSurf = (unsigned char
    *)malloc(bufferSize);
    // Initiaize the buffer
    for(int y = 0; y < bufferSize; y++)
        pSurf[y] = 0;

    //Init texture
    GL_SAFE_CALL(glGenTextures(1, &tex));
    GL_SAFE_CALL(glBindTexture(GL_TEXTURE_2D, tex));
    GL_SAFE_CALL(glTexImage2D(GL_TEXTURE_2D, 0, GL_RGBA, WIDTH, HEIGHT, 0, GL_RGBA,
    GL_UNSIGNED_BYTE, pSurf));
    EGLDisplay eglDisplayHandle = eglGetCurrentDisplay();
    EGLContext eglCtx = eglGetCurrentContext();

    cudaEvent_t cuda_event;
    cudaEventCreateWithFlags(cuda_event, cudaEventDisableTiming);
    EGLAttribKHR eglattrib[] = { EGL_CUDA_EVENT_HANDLE_NV, (EGLAttrib) cuda_event,
    EGL_NONE};
    cudaStream_t* stream;
    cudaStreamCreateWithFlags(&stream, cudaStreamDefault);

    EGLSyncKHR eglsync1, eglsync2;
    cudaEvent_t egl_event;

    // Create the EGL_Image
    EGLint eglImgAttrs[] = { EGL_IMAGE_PRESERVED_KHR, EGL_FALSE, EGL_NONE,
    EGL_NONE };
    EGLImageKHR eglImage = eglCreateImageKHR(eglDisplayHandle, eglCtx,
    EGL_GL_TEXTURE_2D_KHR, (EGLClientBuffer)(intptr_t)tex, eglImgAttrs);

    glTexSubImage2D(GL_TEXTURE_2D, 0, 0, 0, WIDTH, HEIGHT, GL_RGBA,
    GL_UNSIGNED_BYTE, pSurf);
    //Creates an EGLSync object from GL Sync object to track
    //finishing of copy.
    eglsync1 = eglCreateSyncKHR(eglDisplayHandle, EGL_SYNC_FENCE_KHR, NULL);

    //Create CUDA event object from EGLSync obejct
    cuEventCreateFromEGLSync(&egl_event, eglsync1, cudaEventDefault);

    //Waiting on GPU to finish GL copy
    cuStreamWaitEvent(stream, egl_event, 0);

    // Register buffer with CUDA
    cudaGraphicsEGLRegisterImage(&pResource, eglImage,
    cudaGraphicsRegisterFlagsNone);
    //Get CUDA array from graphics resource object
    cudaGraphicsSubResourceGetMappedArray( &pArray, pResource, 0, 0);
    .
    .
    //Create a CUDA surface object from pArray
    struct cudaResourceDesc resDesc;
    memset(&resDesc, 0, sizeof(resDesc));
    resDesc.resType = cudaResourceTypeArray; resDesc.res.array.array = pArray;
    cudaSurfaceObject_t inputSurfObj = 0;
    cudaCreateSurfaceObject(&inputSurfObj, &resDesc);

    dim3 blockSize(32,32);
    dim3 gridSize(width/blockSize.x,height/blockSize.y);
    // Modifies the CUDA array using CUDA surface object
    changeTexture<<<gridSize, blockSize>>>(inputSurfObj, width, height);
    cuEventRecord(cuda_event, stream);
```


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

© 2019 NVIDIA Corporation. All rights reserved.