



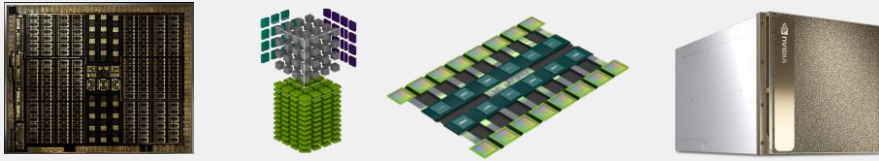
HIGHLIGHTS OF CUDA 10 FOR SUMMIT

March 27, 2019 | OLCF User Conference Call
Steve Abbott

INTRODUCING CUDA 10.0

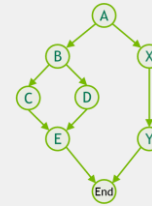
TURING AND NEW SYSTEMS

New GPU Architecture, Tensor Cores, NVSwitch Fabric



CUDA PLATFORM

CUDA Graphs, Vulkan & DX12 Interop, Warp Matrix

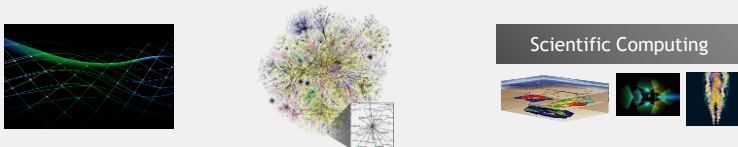


$$D = \begin{matrix} \text{FP16 or FP32} \\ \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix} \\ \text{FP16} \quad \text{FP16} \quad \text{FP16 or FP32} \end{matrix}$$

$D = AB + C$

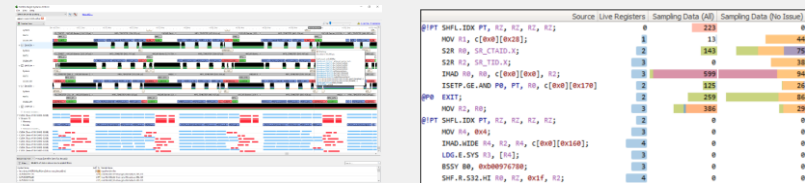
LIBRARIES

GPU-accelerated hybrid JPEG decoding,
Symmetric Eigenvalue Solvers, FFT Scaling



DEVELOPER TOOLS

New Nsight Products - Nsight Systems and Nsight Compute



CUDA DEVELOPMENT ECOSYSTEM

GPU Users

Domain Specialists

Problem Specialists

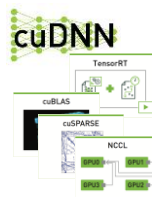
New Algorithm Developers and Optimization Experts



Applications



Frameworks



Libraries



Directives and Standard Languages

CUDA-C++
CUDA Fortran



Extended Standard Languages

Ease of use

Specialized Performance

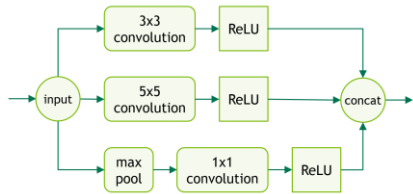
CUDA: Programming Model, GPU Architecture, System Architecture

The background features a complex network of thin, glowing green lines connecting various nodes. The nodes are represented by small, bright green circles of varying sizes and opacities. Some nodes are more prominent, while others are fainter. The lines crisscross the frame, creating a sense of interconnectedness and data flow. The overall aesthetic is futuristic and technical, typical of a digital or network-themed presentation.

PROGRAMMING MODEL

NEW PROGRAMMING MODEL FEATURES

Execution



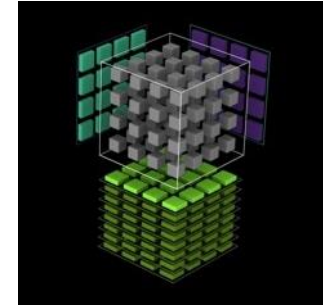
Asynchronous
Task Graphs

Interop



Lightweight Graphics
Interop

Turing



Multi-Precision
Tensor Cores

Precision

```
atomicAdd(&h, (half)1.15f);
half2 hvec(0.94f, -2.13f);
atomicAdd(&h2, hvec);
```

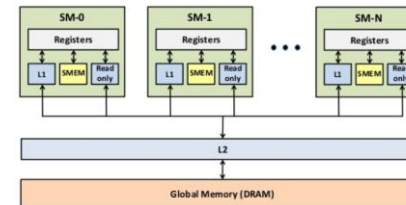
IEEE-754.2008 FP16 Specification

0
0
1
1
1
1
0
0
1
1
0
1
0
1
0
0
0
0
0
 = 0.707031

sign bit exponent (5 bits) mantissa (10 bits)

FP16 Operations

Efficiency



NVCC Enhancements

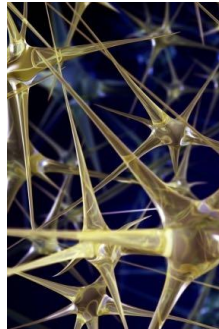
ASYNCHRONOUS TASK GRAPHS

Execution Optimization When Workflow is Known Up-Front

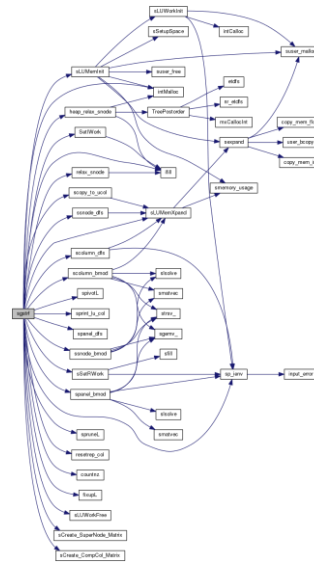
```
// Basic function to test primality.
bool isPrime( size_t n)
{
    if (n == 2) return true;
    if ((n == 1) || ((n % 2) == 0)) return false;
    size_t iters = (unsigned int)sqrt((double)n);
    for (size_t i = 3; i <= iters; i+=2) if (n % i == 0) return false;
    return true;
}

// Compute primes from 1 to 100,000,000.
size_t computePrimes()
{
    size_t primes = 0;
    for ( size_t start = 1; start <= 100000000; ++start)
    {
        if ( isPrime( start) )
        {
            ++primes;
        }
    }
    return primes;
}
```

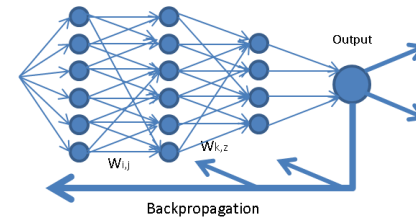
Loop & Function offload



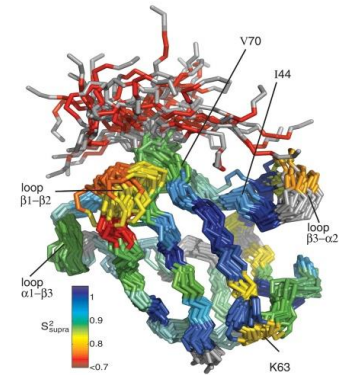
DL Inference



Linear Algebra



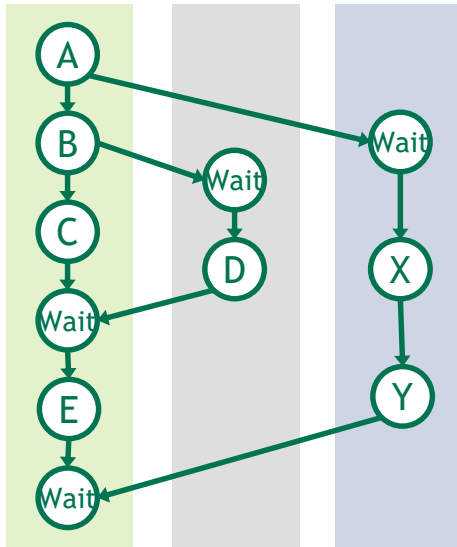
Deep Neural Network Training



HPC Simulation

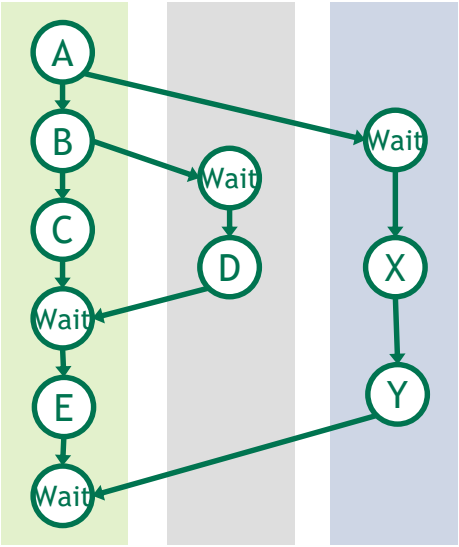
ALL CUDA WORK FORMS A GRAPH

CUDA Work in Streams



ALL CUDA WORK FORMS A GRAPH

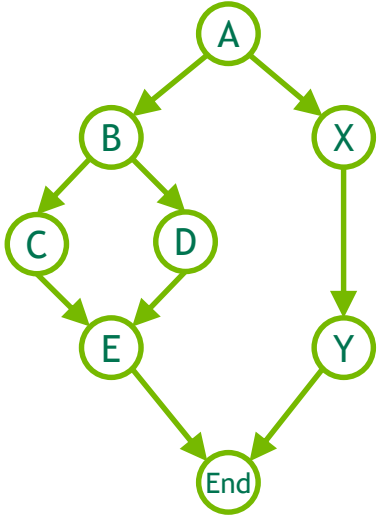
CUDA Work in Streams



Any CUDA stream can be mapped to a graph



Graph of Dependencies



DEFINITION OF A CUDA GRAPH

Graph Nodes Are Not Just Kernel Launches

Sequence of operations, connected by dependencies.

Operations are one of:

Kernel Launch

CUDA kernel running on GPU

CPU Function Call

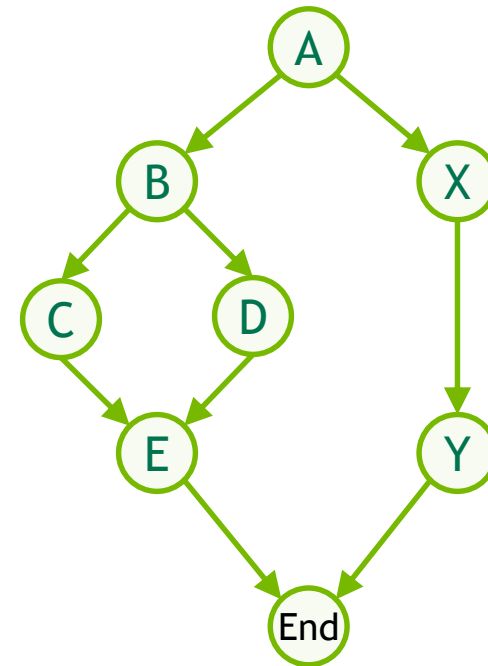
Callback function on CPU

Memcpy/Memset

GPU data management

Sub-Graph

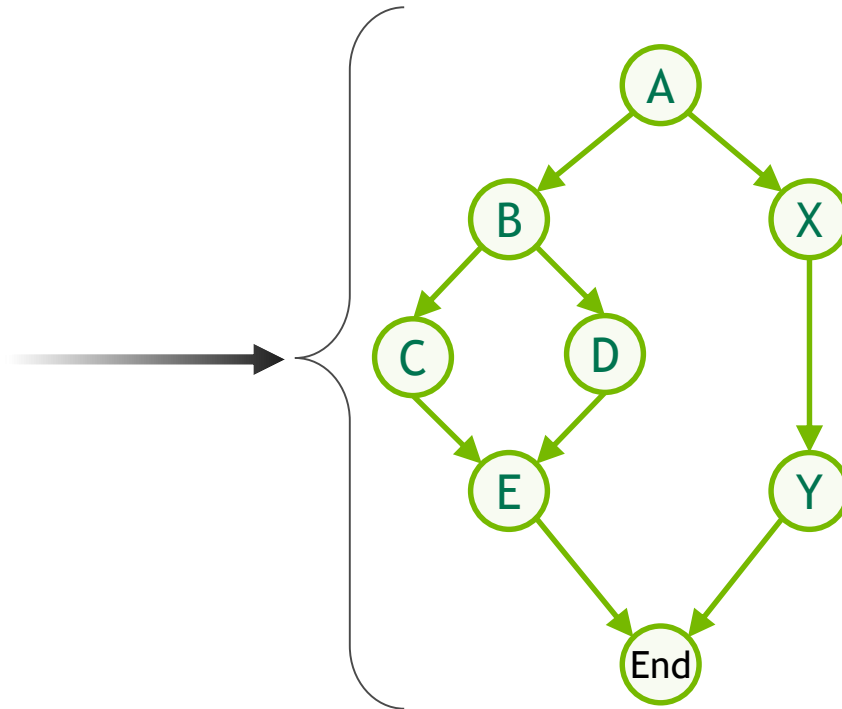
Graphs are hierarchical



NEW EXECUTION MECHANISM

Graphs Can Be Generated Once Then Launched Repeatedly

```
for(int i=0; i<1000; i++) {  
    launch_graph( G );  
}
```

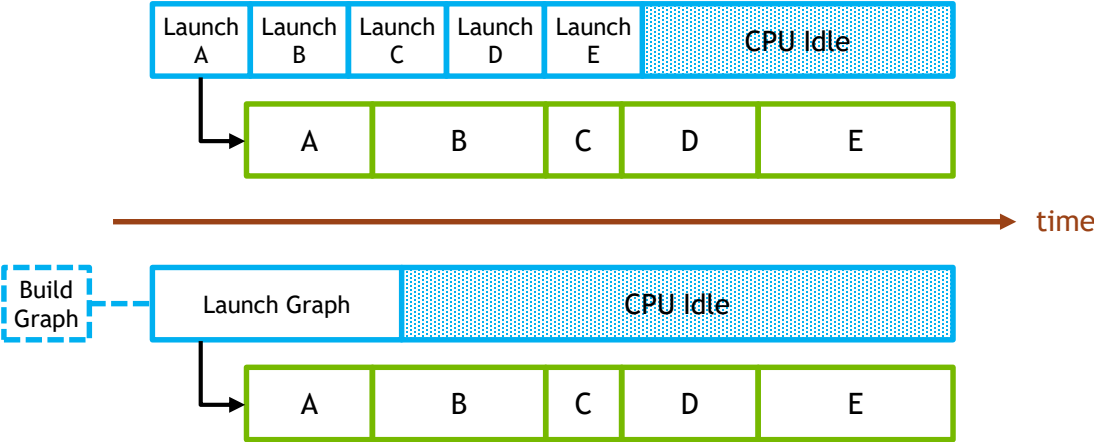


EXECUTION OPTIMIZATIONS

Latency & Overhead Reductions

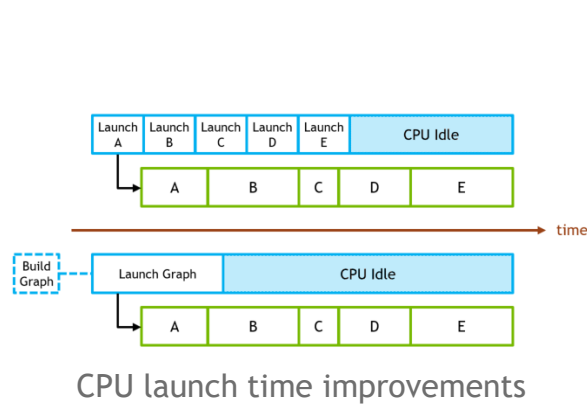
Launch latencies:

- CUDA 10.0 takes at least 2.2us CPU time to launch **each** CUDA kernel on Linux
- Pre-defined graph allows launch of **any number** of kernels in **one single operation**



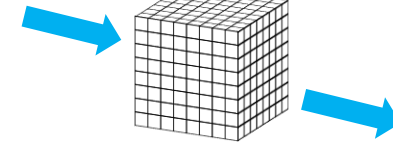
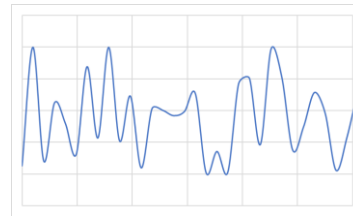
PERFORMANCE IMPACT

Optimizations for Short-Runtime Operations

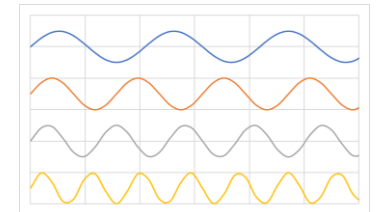


CPU launch time improvements

Typical: **33% faster** than stream launch



Example: Small 3D FFT



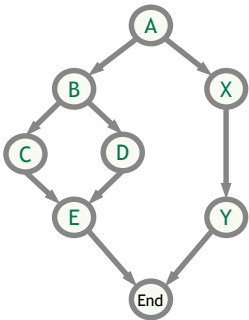
25% end-to-end improvement for 32^3 3D-FFT
(**16us** with stream launch, **12us** with graph launch)

NOTE: Performance impact is workload-dependent

Benefits especially short-running kernels, where overheads account for more runtime

THREE-STAGE EXECUTION MODEL

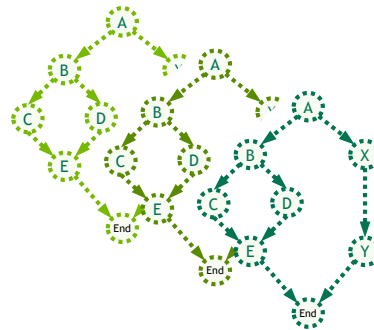
Define



Single Graph “Template”

Created in host code
or built up from libraries

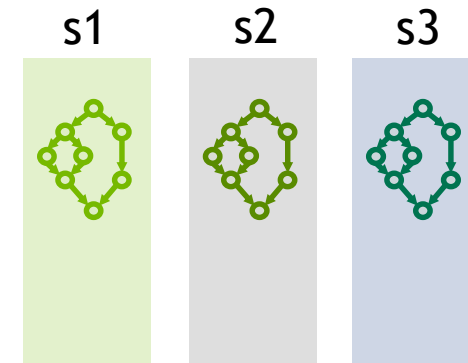
Instantiate



Multiple “Executable Graphs”

Snapshot of template
Sets up & initializes GPU
execution structures
(create once, run many times)

Execute



Executable Graphs
Running in CUDA Streams

Concurrency in graph
is not limited by stream

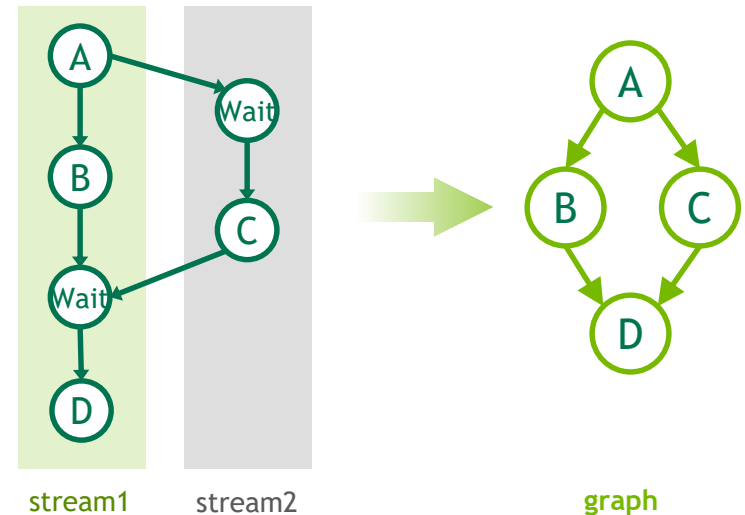
CONVERT CUDA STREAM INTO A GRAPH

Construct a graph from normal CUDA stream syntax

```
// Start by initiating stream capture
cudaStreamBeginCapture(&stream1, cudaStreamCaptureModeGlobal);

// Build stream work as usual
A<<< ..., stream1 >>>();
cudaEventRecord(e1, stream1);
B<<< ..., stream1 >>>();
cudaStreamWaitEvent(stream2, e1);
C<<< ..., stream2 >>>();
cudaEventRecord(e2, stream2);
cudaStreamWaitEvent(stream1, e2);
D<<< ..., stream1 >>>();

// Now convert the stream to a graph
cudaStreamEndCapture(stream1, &graph);
```



CONVERT CUDA STREAM INTO A GRAPH

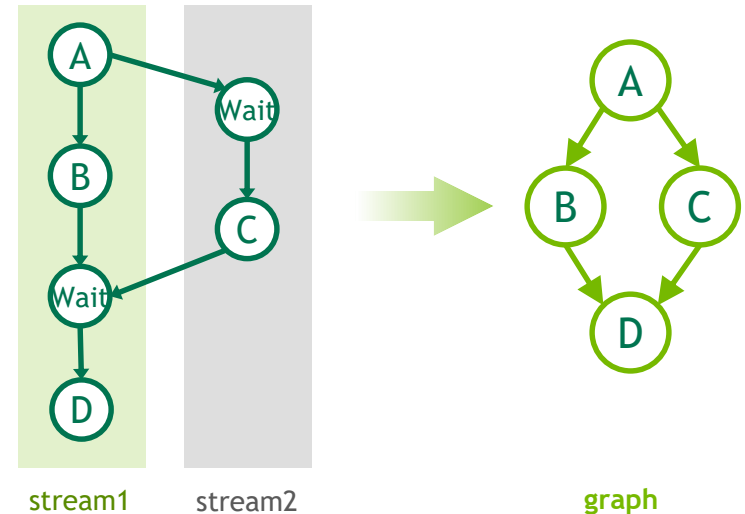
Construct a graph from normal CUDA stream syntax

```
// Start by initiating stream capture
cudaStreamBeginCapture(&stream1, cudaStreamCaptureModeGlobal);

// Build stream work as usual
A<<< ..., stream1 >>>();
cudaEventRecord(e1, stream1);
B<<< ..., stream1 >>>();
cudaStreamWaitEvent(stream2, e1);
C<<< ..., stream2 >>>();
cudaEventRecord(e2, stream2);
cudaStreamWaitEvent(stream1, e2);
D<<< ..., stream1 >>>();

// Now convert the stream to a graph
cudaStreamEndCapture(stream1, &graph);
```

Capture follows inter-stream dependencies to create forks & joins



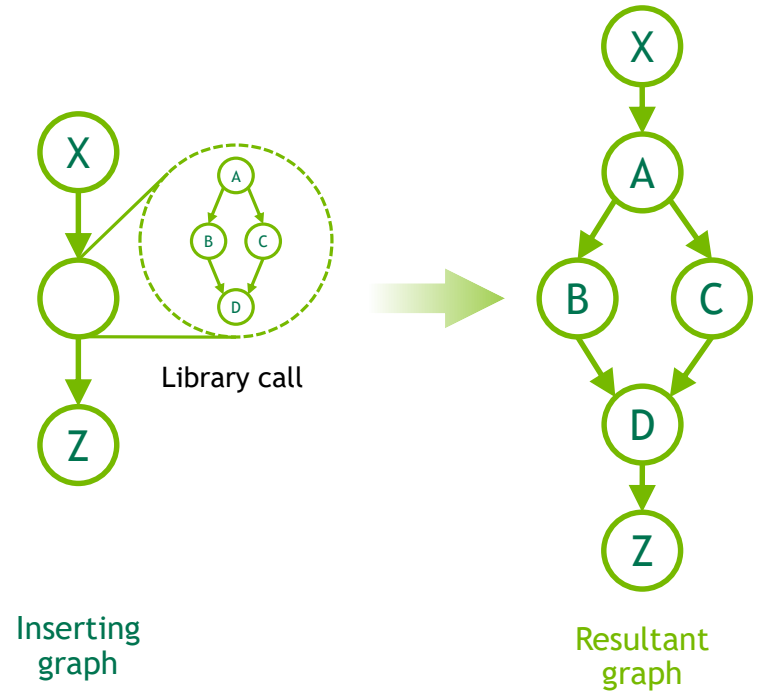
CAPTURE EXTERNAL WORK

Stream Capture Continues Into Library Calls

```
// Start by initiating stream capture
cudaStreamBeginCapture(&stream, cudaStreamCaptureModeGlobal);

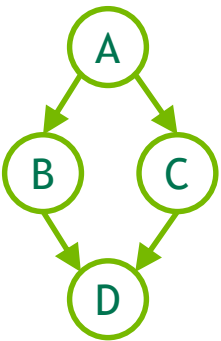
// Captures my kernel launches, recurse into library calls
X<<< ..., stream >>>();
libraryCall(stream);           // Launches A, B, C, D
Z<<< ..., stream >>>();

// Now convert the stream to a graph
cudaStreamEndCapture(stream, &graph);
```



CREATE GRAPHS DIRECTLY

Map Graph-Based Workflows Directly Into CUDA



Graph from
framework



```
// Define graph of work + dependencies
cudaGraphCreate(&graph);

cudaGraphAddNode(graph, kernel_a, {}, ...);
cudaGraphAddNode(graph, kernel_b, { kernel_a }, ...);
cudaGraphAddNode(graph, kernel_c, { kernel_a }, ...);
cudaGraphAddNode(graph, kernel_d, { kernel_b, kernel_c }, ...);

// Instantiate graph and apply optimizations
cudaGraphInstantiate(&instance, graph);

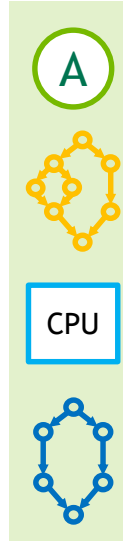
// Launch executable graph 100 times
for(int i=0; i<100; i++)
    cudaGraphLaunch(instance, stream);
```

GRAPH EXECUTION SEMANTICS

Order Graph Work With Other Non-Graph CUDA Work

```
launchWork(cudaGraphExec_t i1, cudaGraphExec_t i2,  
           CPU_Func cpu, cudaStream_t stream) {  
  
    A <<< 256, 256, 0, stream >>>();           // Kernel launch  
    cudaGraphLaunch(i1, stream);                // Graph1 launch  
    cudaStreamAddCallback(stream, cpu);        // CPU callback  
    cudaGraphLaunch(i2, stream);                // Graph2 launch  
  
    cudaStreamSynchronize(stream);  
}
```

stream



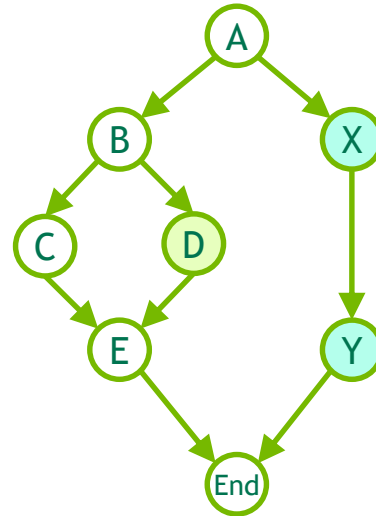
If you can put it in a CUDA stream, you can run it together with a graph

GRAPHS IGNORE STREAM SERIALIZATION RULES

Launch Stream Is Used Only For Ordering With Other Work



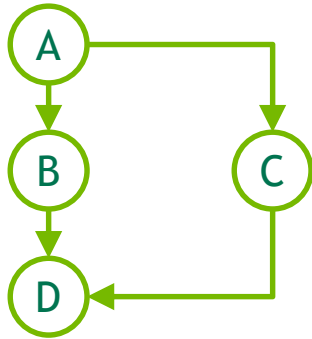
Branches in graph still execute concurrently even though graph is launched into a stream



CROSS-DEVICE DEPENDENCIES

Graphs May Span Multiple GPUs

Multi-Device Execution



GPU 0



GPU 1

Heterogeneous Execution



CUDA is closest to the O/S and the hardware

- Can optimize **multi-device** dependencies
- Can optimize **heterogeneous** dependencies
- Define locality per-node

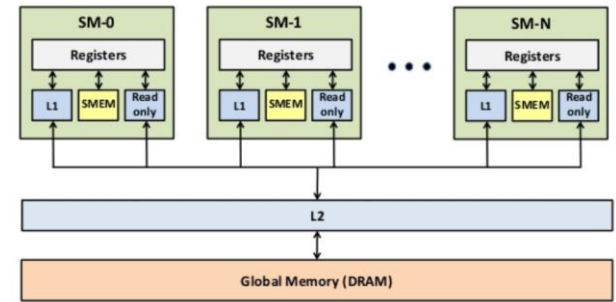
NVCC IN CUDA 10

Improving Efficiency

Extensible Whole Program (-ewp) mode
compilation support

Enables efficient compilation with use of CUDA
run-time device library

Compiler optimization and code generation
heuristics tuning for Volta and Turing



Efficient Code Generation for
Chip Architecture

ENHANCED HALF-PRECISION FUNCTIONALITY

Includes Limited *half* Type Support For CPU Code

Half-precision atomic ADD
(Volta+) (round-to-nearest mode)

```
half atomicAdd(half *address, half val);  
half2 atomicAdd(half2 *address, half2 val);
```

Host-side conversion operators
between *float* and *half* types

```
half pi = 3.1415f;           // Convert float to half  
float fPI = (float)hPI;     // Convert half to float
```

Host-side construction and
assignment operators for
half and *half2* types

```
half pi = 3.1415f;  
half also_pi = pi;          // Assign half to half  
half2 vector_pi(pi, also_pi); // Construct half2 from half
```

NOTE: Half-precision *arithmetic* operations remain only available in device code



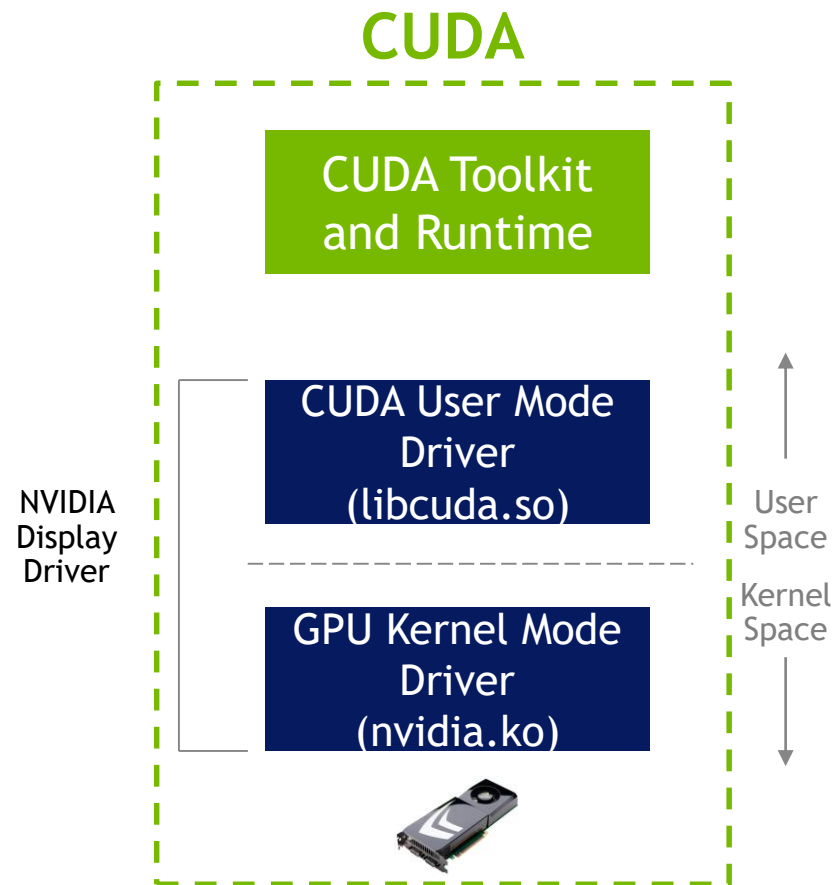
CUDA DEPLOYMENT

CUDA INSTALLED COMPONENTS

CUDA is comprised of three components

1. CUDA Toolkit (build applications)
2. CUDA User Mode Driver (run applications)
3. NVIDIA Kernel Mode Driver (run applications)

Note that components 2 and 3 are delivered together in the NVIDIA Display Driver package

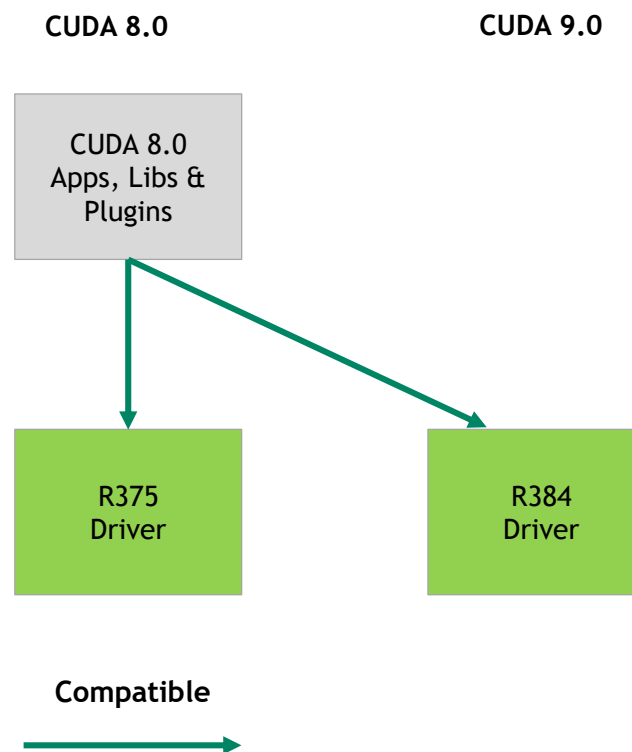


CUDA COMPATIBILITY - TODAY

Older CUDA Version Runs on Newer Display Driver

CUDA driver API is backward compatible but not forward compatible

- ▶ Each CUDA release has a minimum driver requirement
- ▶ Applications compiled against a particular version of CUDA API will work on later driver releases
- ▶ E.g.
 - ▶ CUDA 8.0 needs \geq R375
 - ▶ CUDA 9.0 needs \geq R384

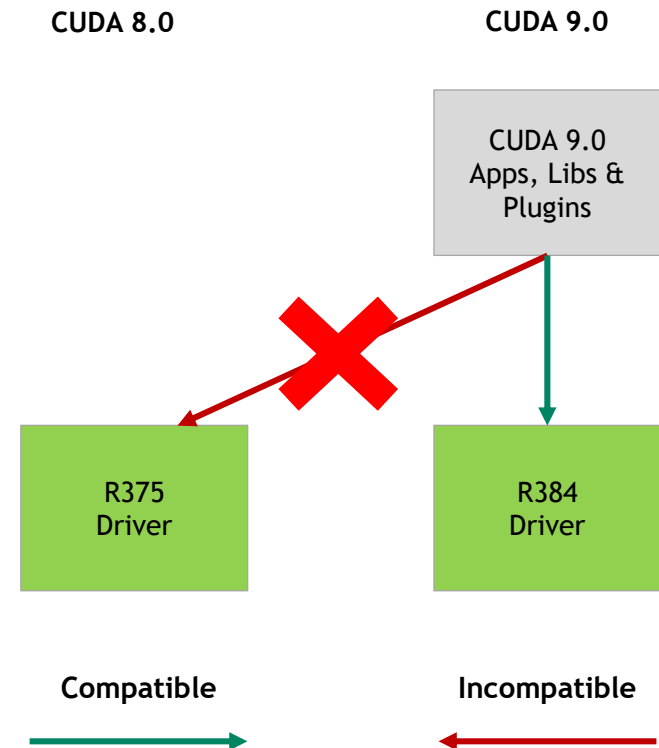


CUDA COMPATIBILITY - TODAY

Newer CUDA Version **DOES NOT** Run on Older Display Driver

CUDA driver API is backward compatible but not forward compatible

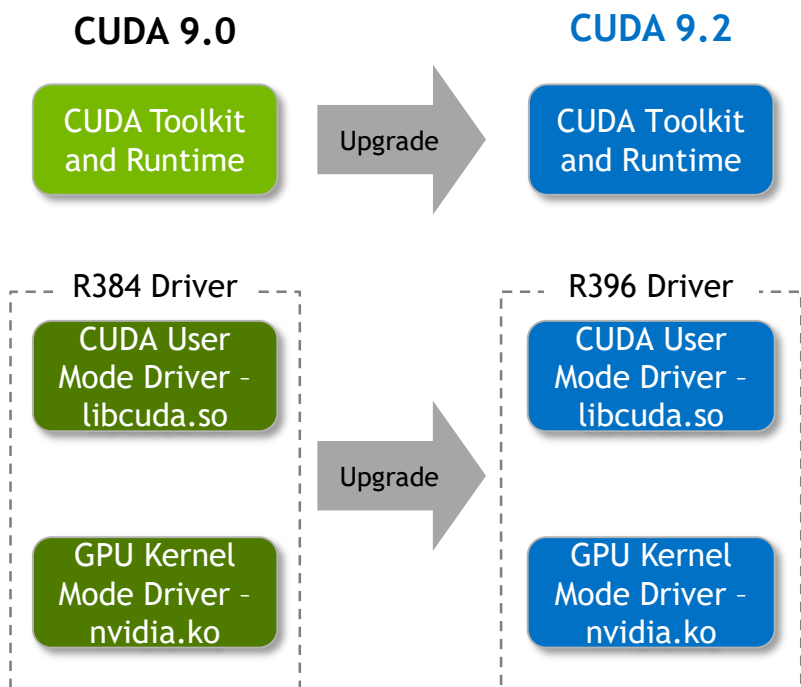
- ▶ Each CUDA release has a minimum driver requirement
- ▶ Applications compiled against a particular version of CUDA API will work on later driver releases
- ▶ E.g.
 - ▶ CUDA 8.0 needs \geq R375
 - ▶ CUDA 9.0 needs \geq R384



CUDA COMPATIBILITY - UPGRADE PATHS

Backward Compatibility

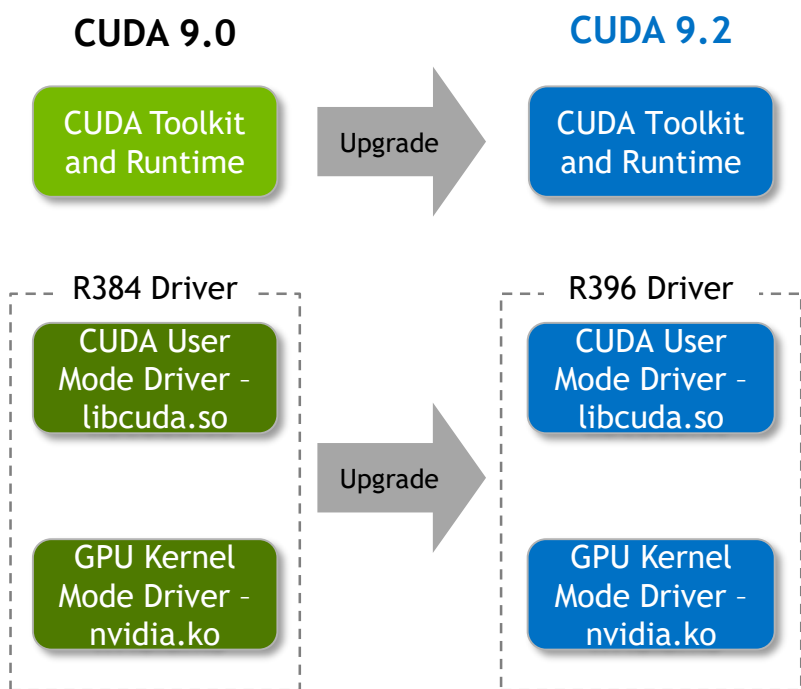
Upgrade *both* toolkit and driver



CUDA COMPATIBILITY - UPGRADE PATHS

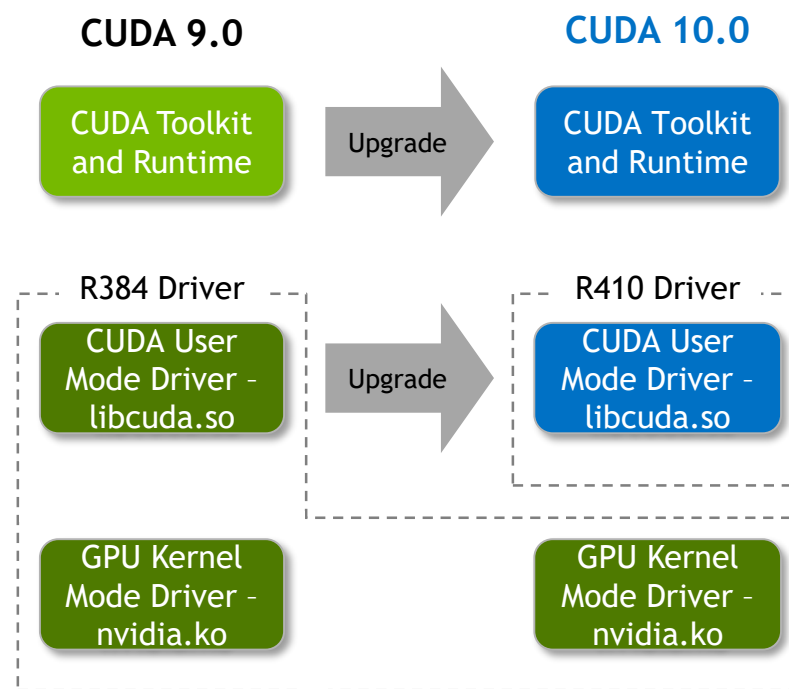
Backward Compatibility

Upgrade *both* toolkit and driver



NEW Forward Compatibility *Option*

Upgrade *only* user-mode CUDA components*



*requires new 'cuda-compat-10-0' package

CUDA COMPATIBILITY - UPGRADE PATHS

Starting with CUDA 10.0

New compatibility platform upgrade path available

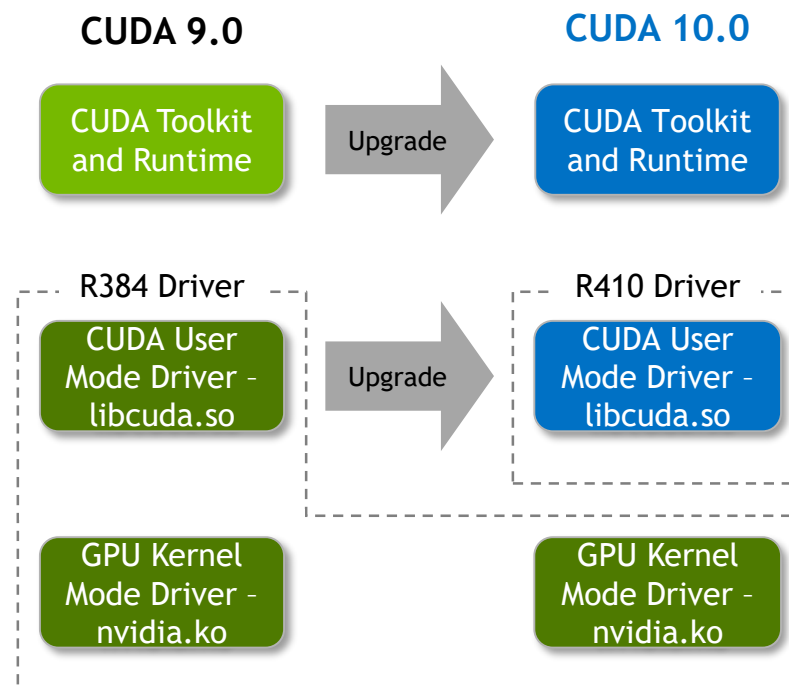
- ▶ Use newer CUDA toolkits on older driver installs
- ▶ Compatibility only with specific older driver versions

System requirements

- ▶ Tesla GPU support only - no Quadro or GeForce
- ▶ Only available on Linux

NEW Forward Compatibility *Option*

Upgrade *only* user-mode CUDA components*



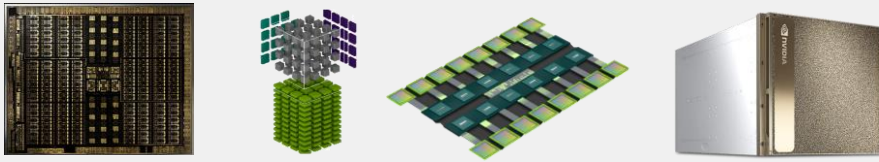
*requires new 'cuda-compat-10-0' package

CUDA 10.1 - COMING TO SUMMIT SOON!

<https://developer.nvidia.com/cuda-toolkit>

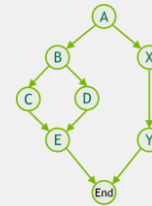
TURING AND NEW SYSTEMS

New GPU Architecture, Tensor Cores, NVSwitch Fabric



CUDA PLATFORM

CUDA Graphs, Vulkan & DX12 Interop, Warp Matrix

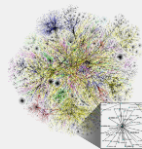
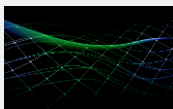


$$D = \begin{matrix} \text{FP16 or FP32} \\ \begin{pmatrix} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{pmatrix} \begin{pmatrix} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{pmatrix} + \begin{pmatrix} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{pmatrix} \\ \text{FP16} & \text{FP16} & \text{FP16 or FP32} \end{matrix}$$

$D = AB + C$

LIBRARIES

GPU-accelerated hybrid JPEG decoding,
Symmetric Eigenvalue Solvers, FFT Scaling



DEVELOPER TOOLS

New Nsight Products - Nsight Systems and Nsight Compute



Source	Line	Registers	Sampling Data (All)	Sampling Data (No Issue)
0x10000000	0	0	223	0
0x10000000	1	0	13	44
0x10000000	2	0	180	206
0x10000000	3	0	3	38
0x10000000	4	0	599	54
0x10000000	5	0	325	24
0x10000000	6	0	288	86
0x10000000	7	0	386	25
0x10000000	8	0	0	0
0x10000000	9	0	0	0
0x10000000	10	0	0	0
0x10000000	11	0	0	0
0x10000000	12	0	0	0
0x10000000	13	0	0	0
0x10000000	14	0	0	0
0x10000000	15	0	0	0
0x10000000	16	0	0	0
0x10000000	17	0	0	0
0x10000000	18	0	0	0
0x10000000	19	0	0	0
0x10000000	20	0	0	0