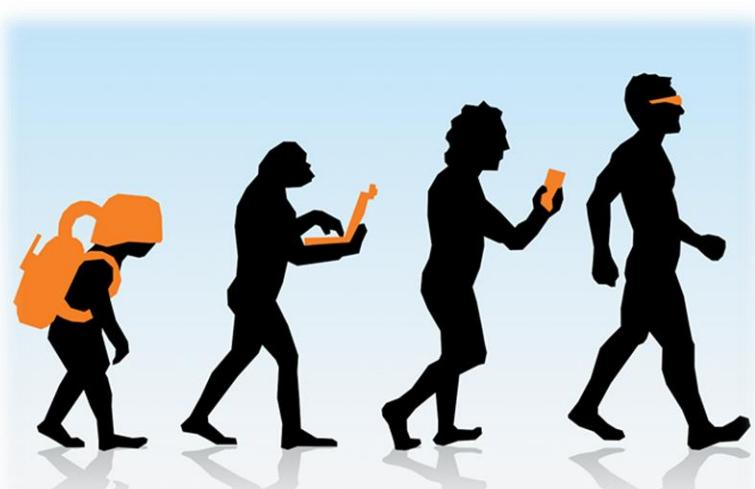


# CUDA Kernels with C++

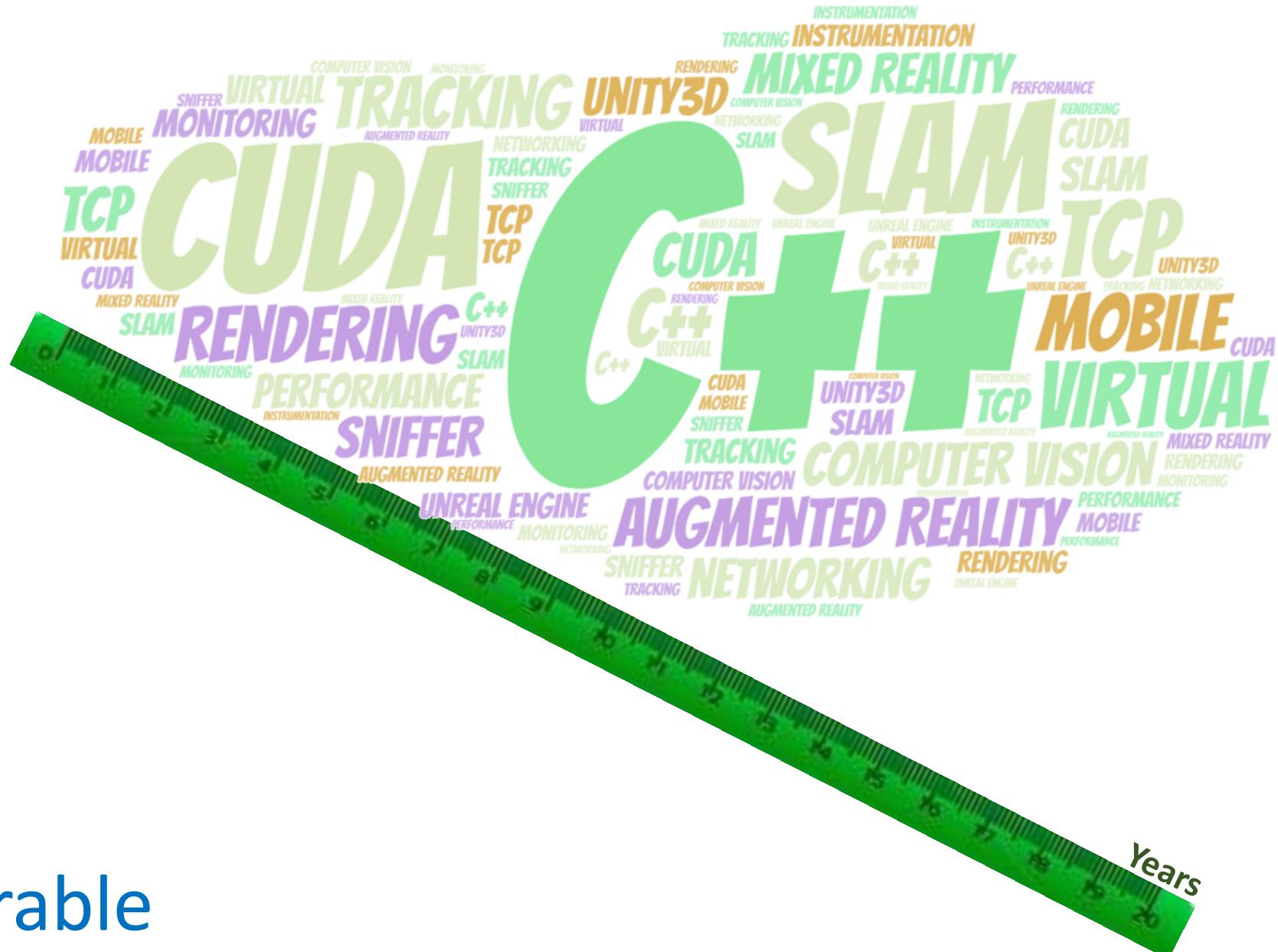
Michael Gopshtein

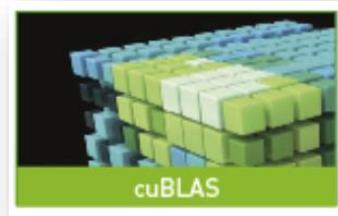
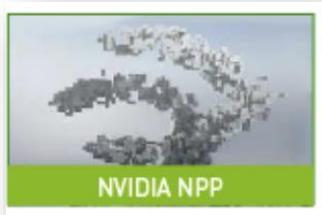
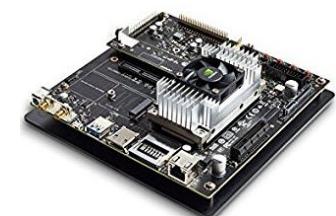
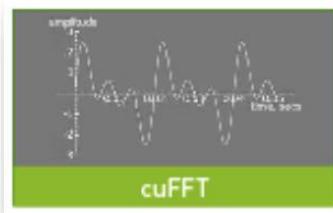
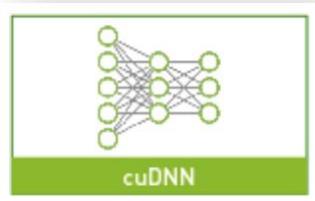
Sep 2018

# About me



Photorealistic wearable  
Augmented Reality experience.

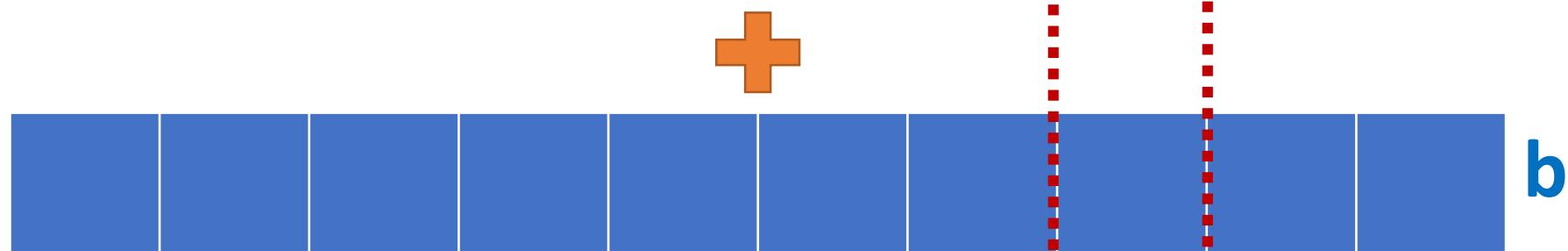




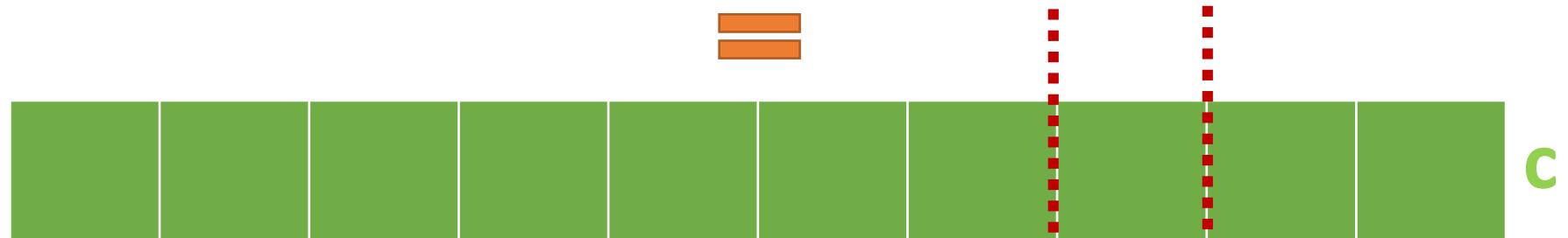
# Setting the Ground



a



b



c

#7 CUDA “Thread”

# Vector Addition in CUDA

```
__global__ void addKernel(int *c, const int *a, const int *b) {  
    int idx = (blockIdx.x * blockDim.x) + threadIdx.x;  
    c[idx] = a[idx] + b[idx];  
}
```

Kernel/Device Code

```
int main() {  
    //...  
    addKernel<<<blocks, 32>>>(dev_c, dev_a, dev_b);  
    //...  
}
```

Host/CPU Code

# Vector Addition in CUDA

```
__global__ void addKernel(int *c, int *a, const int *b) {  
    int idx = (blockIdx.x * blockDim.x) + threadIdx.x;  
    c[idx] = a[idx] + b[idx];  
}  
  
int main() {  
    //...  
    addKernel<<<blocks, 32>>>(dev_c, dev_a, dev_b);  
    //...  
}
```

GPU Accelerated Computing with C and C++

This is C  
Where are the pluses?

Professional  
**CUDA® C**  
Programming



# Different Ways to Use C++



CUDA C/C++



OpenACC  
OpenMP

*Other  
libraries*

# Vector Addition in CUDA

```
__global__ void addKernel(int *c, const int *a, const int *b) {  
    int idx = (blockIdx.x * blockDim.x) + threadIdx.x;  
    c[idx] = a[idx] + b[idx];  
}  
  
int main() {  
    //...  
    addKernel<<<blocks, 32>>>(dev_c, dev_a, dev_b);  
    //...  
}
```

What if we have  
*float* arrays?

# C Way

```
__global__ void addKernel(int *c, const int *a, const int *b) {  
    int idx = (blockIdx.x * blockDim.x) + threadIdx.x;  
    c[idx] = a[idx] + b[idx];  
}  
  
__global__ void addKernelF(float *c, const float *a, const float *b) {  
    int idx = (blockIdx.x * blockDim.x) + threadIdx.x;  
    c[idx] = a[idx] + b[idx];  
}
```

# In C++ it's easy!

```
Device
template<typename T>
__global__ void addKernel(T *c, const T *a, const T *b) {
    int idx = (blockIdx.x * blockDim.x) + threadIdx.x;
    c[idx] = a[idx] + b[idx];
}
```

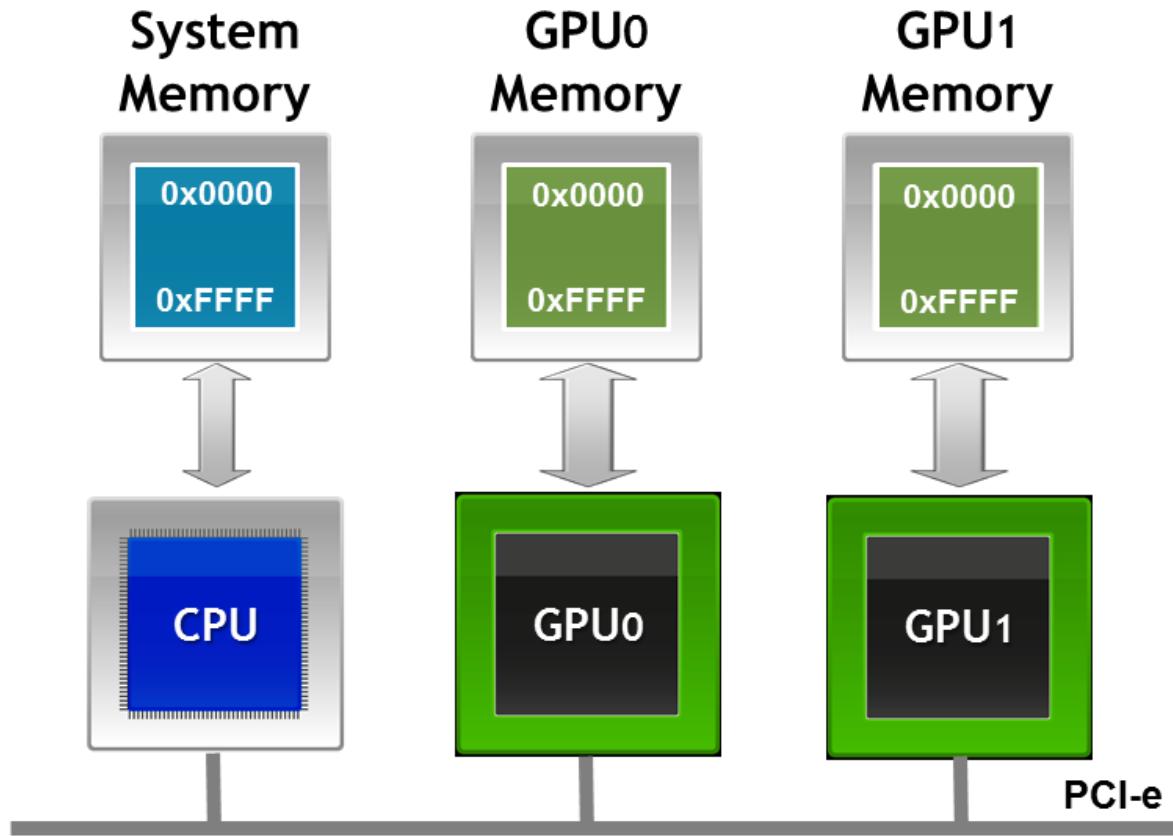
```
Host
addKernel<int><<<blocks, 32>>>(dev_c, dev_a, dev_b);
```

```
addKernel<<<blocks, 32>>>(dev_c, dev_a, dev_b);
```

```
template<typename T>
__global__ void addKernel(T *c, const T *a, const T *b) {
    int idx = (blockIdx.x * blockDim.x + threadIdx.x;
    c[idx] = a[idx] + b[idx];
}
```

What memory  
does it point to?

```
addKernel<<<blocks, 32>>>(dev_c, dev_a, dev_b);
```



***cudaMalloc*** allocates memory on the GPU

***cudaMemcpy*** copies the vectors to/from GPU

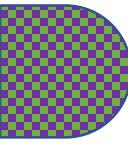
# Compiles, but fails in runtime

```
Device  
template<typename T>  
__global__ void addKernel(T *c, const T *a, const T *b) {...}  
  
Host  
int main() {  
    const int a[SIZE] = {1, 2, ...};  
    const int b[SIZE] = {10, 20, ...};  
    int c[SIZE];  
  
    addKernel<<<blocks, 32>>>(c, a, b);  
}
```



# Let's use explicit device memory pointers

```
template<typename T>
__global__ void addKernel(
    DevicePtr<T> c,
    DevicePtr<const T> a,
    DevicePtr<const T> b
){...}
```



```
template<typename T>
class DevicePtr {
    T *_p = nullptr;

public:
    __device__ __host__ __inline__ explicit DevicePtr(T *p) : _p(p) {}

    //...
};
```

Explicit creation  
from raw T\*



```
template<typename T>
__host__ inline auto MakeDevicePtr(T* p) {
    return DevicePtr<T>(p);
}
```

Convenience  
global function

# Simple usage

```
void process(const int *a, const int *b, int *c) {  
    int *aDev;  
    cudaMalloc(&aDev, LEN);  
    cudaMemcpy(aDev, a, LEN, cudaMemcpyHostToDevice);  
    //... same for bDev(alloc+copy) and cDev(alloc)  
  
    addKernel<<<blocks, 32>>>(MakeDevicePtr(cDev),  
        MakeDevicePtr(aDev), MakeDevicePtr(bDev));  
  
    cudaMemcpy(c, cDev, LEN, cudaMemcpyDeviceToHost);  
    cudaFree(aDev); // free bDev, cDev  
}
```

# Even simpler usage

```
void process(const int *a, const int *b, int *c) {  
    auto aDev = DeviceMemory<int>::AllocateElements(NUM);  
    CopyElements(aDev, a, NUM);  
    //... same for bDev(alloc+copy) and cDev(alloc)  
  
    addKernel<<<blocks, 32>>>(cDev, aDev, bDev);  
  
    CopyElements(c, cDev, LEN);  
}
```

```
template<typename T>
class DeviceMemory {
    T *_p = nullptr;
    DeviceMemory(std::size_t bytes) { cudaMalloc(&_p, _bytes); }

public:
    static DeviceMemory AllocateElements(std::size_t n) {return {n*sizeof(T)}; }
    static DeviceMemory AllocateBytes(std::size_t bytes) {return {bytes}; }
    ~DeviceMemory() { if (_p) {cudaFree(_p);} }

    operator DevicePtr<T>() const {
        return DevicePtr<T>(_p);
    }
};
```

`DevicePtr<int> a;`

`DevicePtr<const int> b = a;` ✓

`DevicePtr<int> c = b;` ✗

`DevicePtr<char> d = a;` ✗

## <type\_traits>

```
template<typename T>
class DevicePtr {
```

```
    T *_p = nullptr;
```

```
template<typename T1,
         typename = std::enable_if_t<std::is_convertible_v<T1*, T*>>>
__device__ __host__ DevicePtr(const DevicePtr<T1> &dp)
    : _p(dp.get())
{}
```

```
};
```

DevicePtr<T1> → DevicePtr<T>

iff

T1\* → T\*

# Static Polymorphism

```
template<typename T>
struct BinaryOpPlus {
    __device__ T operator()(T t1, T t2) const override { return t1 + t2; }
};

template<template<typename> typename OP, typename T> __global__
void addKernel(DevPtr<T> c, DevPtr<const T> a, DevPtr<const T> b) {
    OP<T> op;
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    c[i] = op(a[i], b[i]);
}

addKernel<BinaryOpPlus><<<blocks, 32>>>(cDev, aDev, bDev);
```

Professional  
**CUDA C**  
Programming



<static>  
Polymorphism

Professional  
**CUDA C**  
Programming

<static>  
**Polymorphism**

**Dynamic  
Polymorphism**

```
template<typename T>
struct BinaryOp {
    virtual __device__ T operator()(T t1, T t2) const = 0;
};

template<typename T>
struct BinaryOpPlus : public BinaryOp<T> {
    __device__ T operator()(T t1, T t2) const override { return t1 + t2; }
};

template<typename T>
__device__ void worker(DevPtr<T> c, DevPtr<const T> a, DevPtr<const T> b, const BinaryOp<T> &op) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    c[i] = op(a[i], b[i]);
}

template<typename T>
__global__ void addKernel(DevPtr<T> c, DevPtr<const T> a, DevPtr<const T> b) {
    worker(c, a, b, BinaryOpPlus<T>{});
}
```

*It is not allowed to pass as an argument to a global function an object of a class derived from virtual base classes*

*but the business logic is in the  
CPU code...*



# Solution 1: from template to virtual

```
template<typename T>
struct BinaryOpPlus : public BinaryOp<T> {

    template<typename T> __device__
    void worker(DevPtr<T> c, DevPtr<const T> a, DevPtr<const T> b,
                const BinaryOp<T> &op);

    template<template<typename> typename OP, typename T> __global__
    void addKernel(DevPtr<T> c, DevPtr<const T> a, DevPtr<const T> b) {
        worker(c, a, b, OP<T>{});
    }

    addKernel<BinaryOpPlus><<<blocks, 32>>>(cDev, aDev, bDev);
```

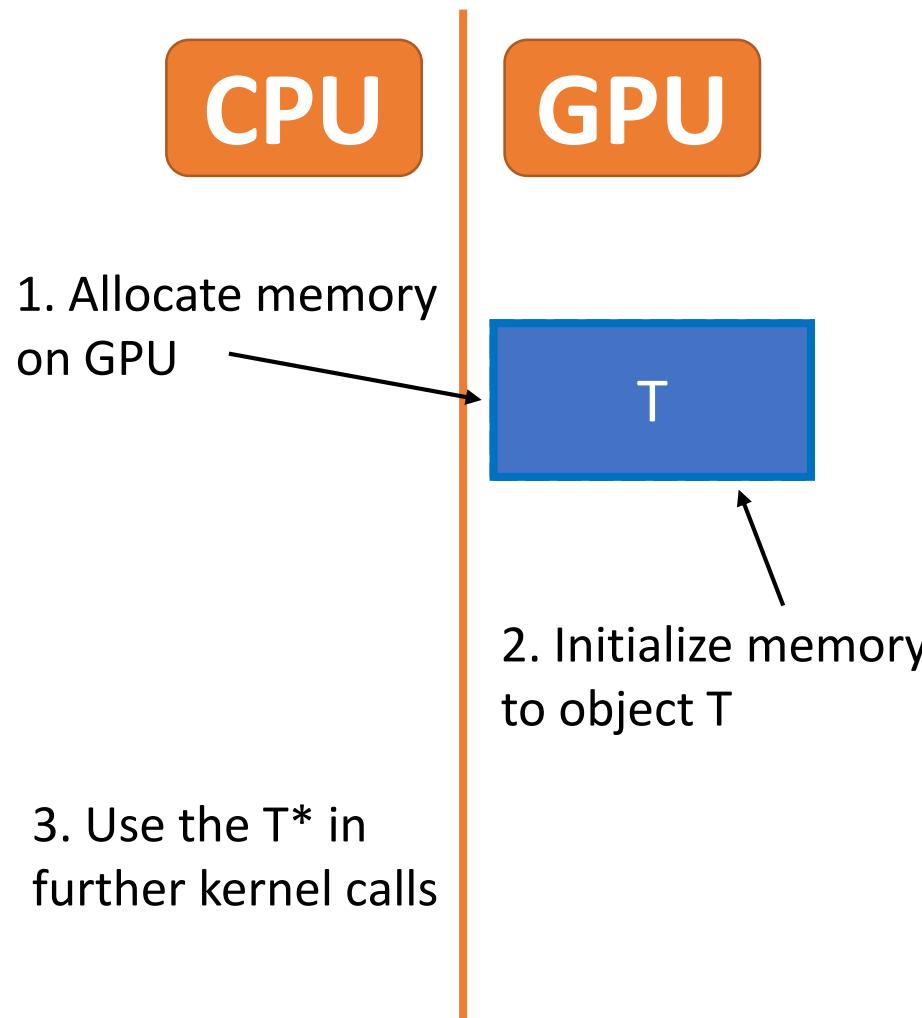
## Solution 2: dynamic allocation

```
template<typename T>
struct BinaryOpPlus : public BinaryOp<T>

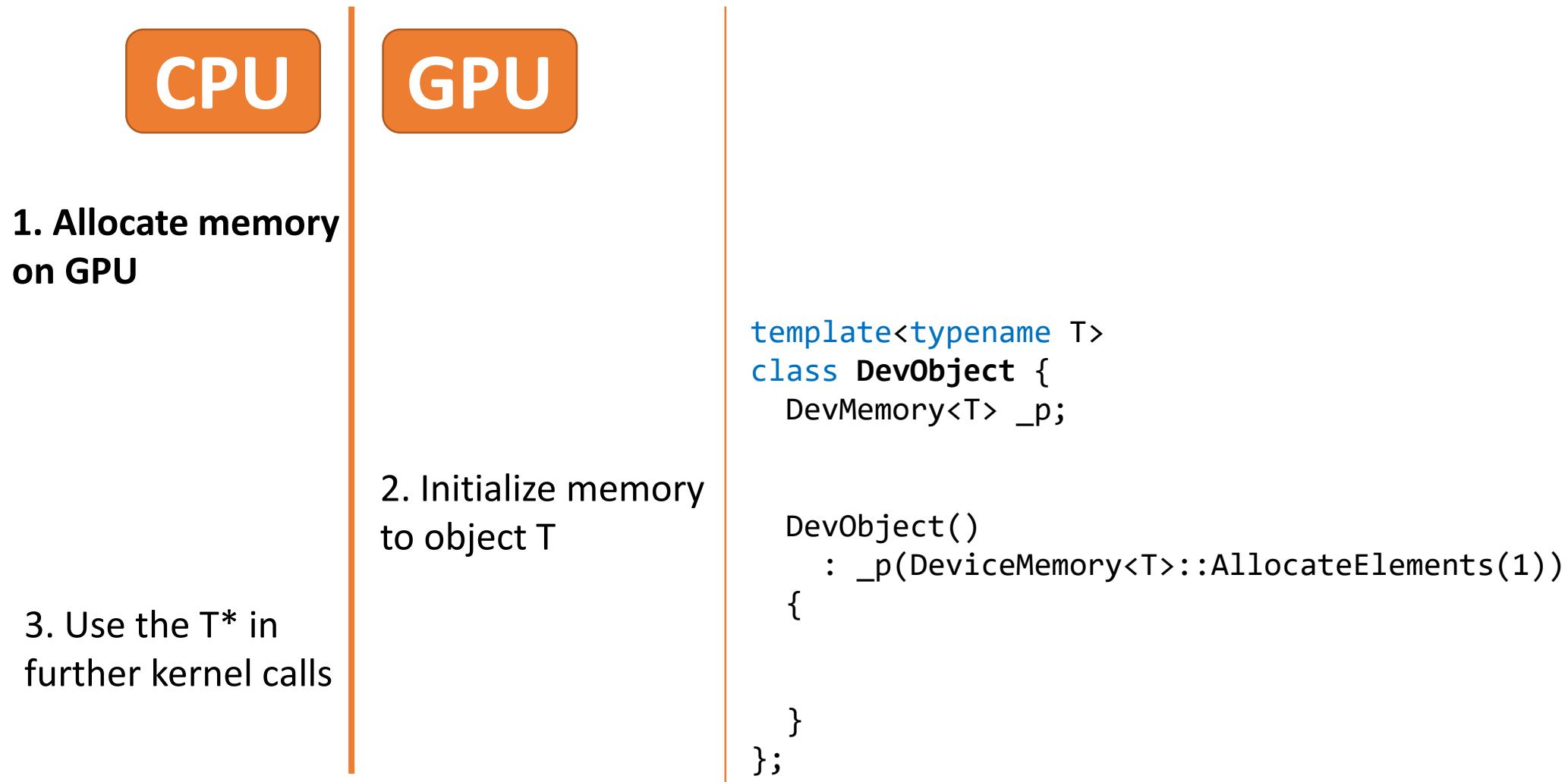
template<typename T> __global__
void addKernel(DevPtr<T> c, DevPtr<const T> a, DevPtr<const T> b,
               DevPtr<const BinaryOp<T>> op)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    c[i] = (*op)(a[i], b[i]);
}

// create BinaryOpPlus on device
addKernel<<<...>>>(cDev, aDev, bDev, op);
```

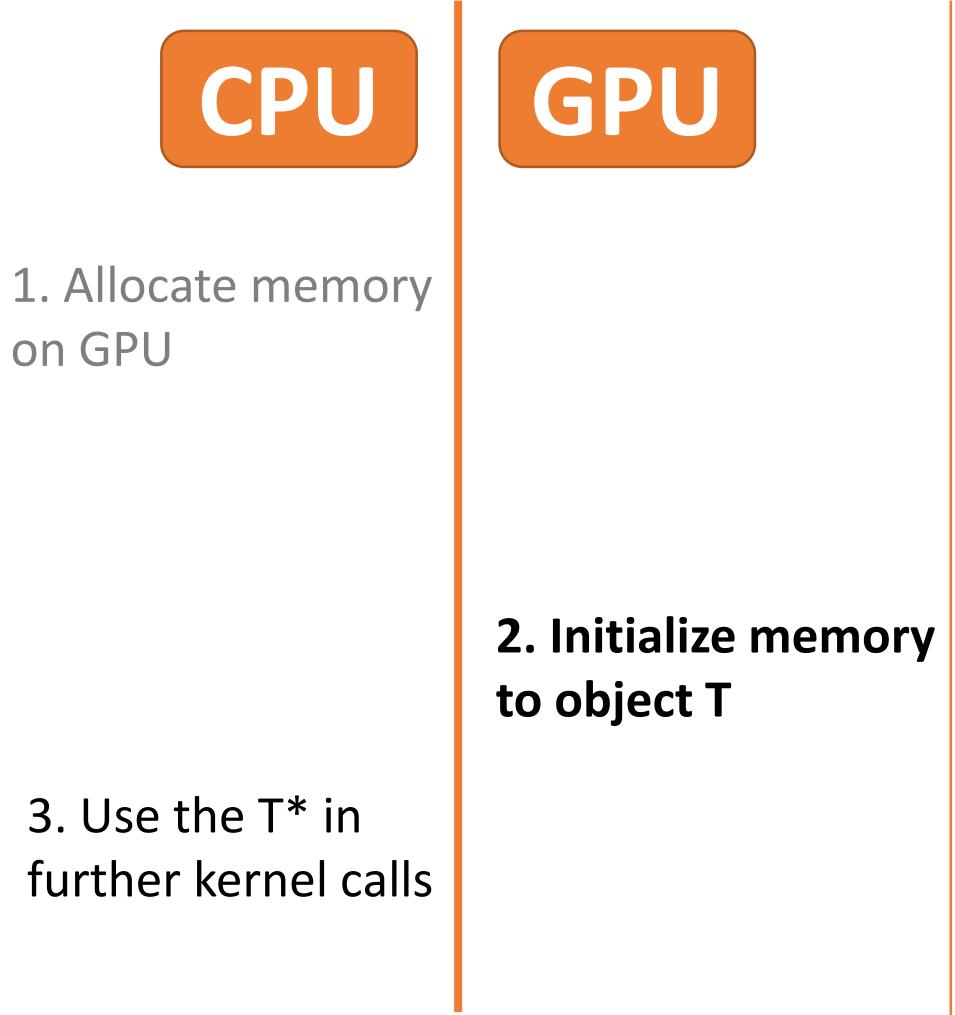
# Dynamic allocation in CUDA



# Dynamic allocation in CUDA



# Dynamic allocation in CUDA

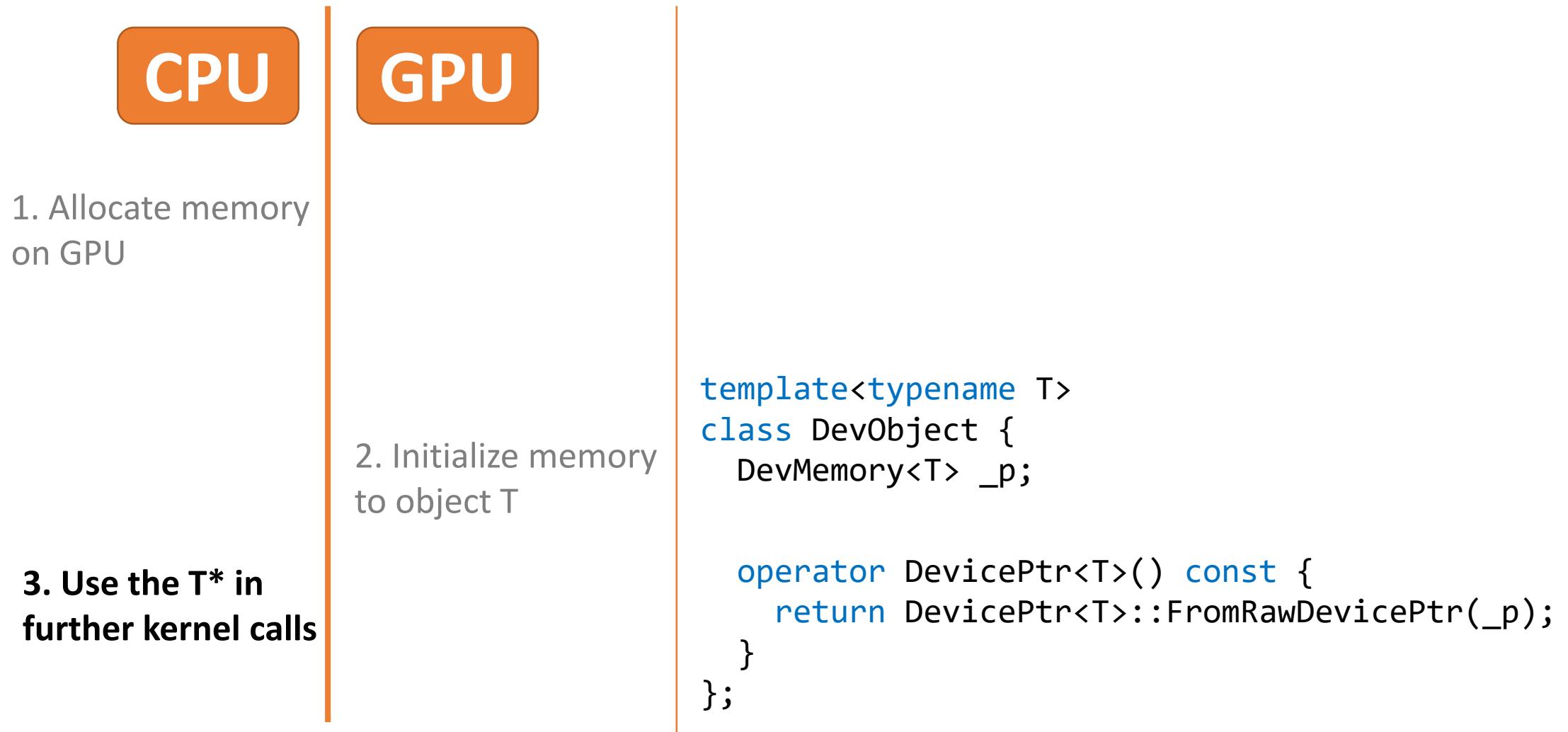


```
namespace detail {
    template<typename T, typename... ARGS> __global__
    void AllocateObject(DevPtr<T> p, ARGS... args) {
        new (p.get()) T(args...);
    }
}

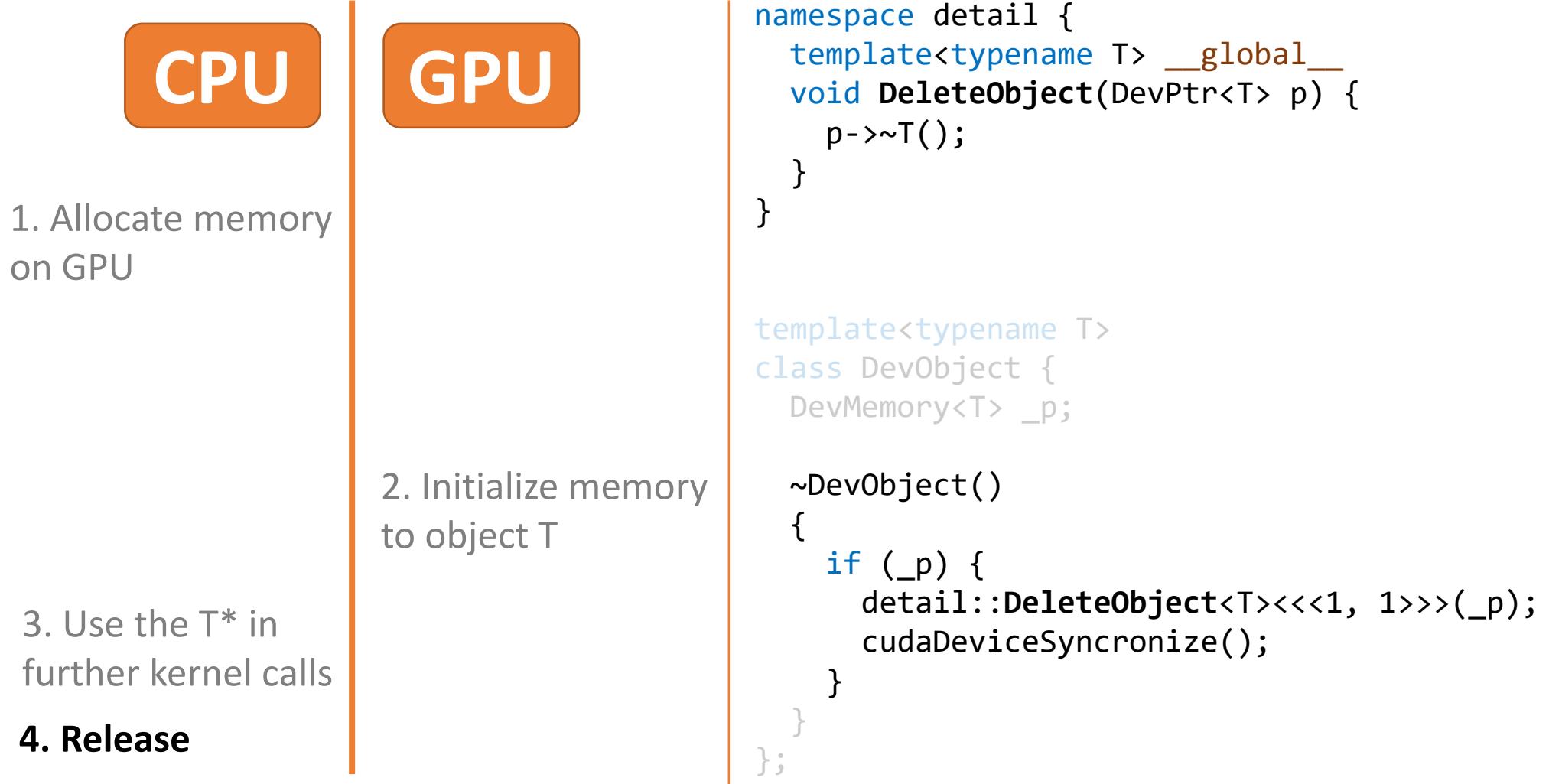
template<typename T>
class DevObject {
    DevMemory<T> _p;

    template<typename... ARGS>
    DevObject(ARGS... args)
        : _p(DeviceMemory<T>::AllocateElements(1))
    {
        detail::AllocateObject<T><<<1, 1>>>(_p, args...);
        cudaDeviceSynchronize();
    }
};
```

# Dynamic allocation in CUDA



# Dynamic allocation in CUDA



Device

Host

<static>  
**Polymorphism**

**Dynamic  
Polymorphism**

**new/delete**

You can just use  
*malloc/free*  
and  
*new/delete*  
in the kernel code

<static>  
**Polymorphism**

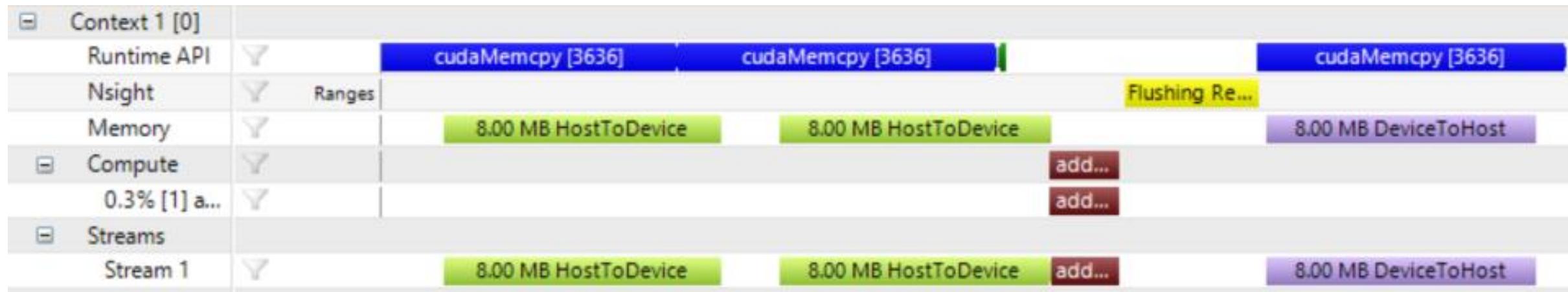
**Dynamic  
Polymorphism**

**new/delete**

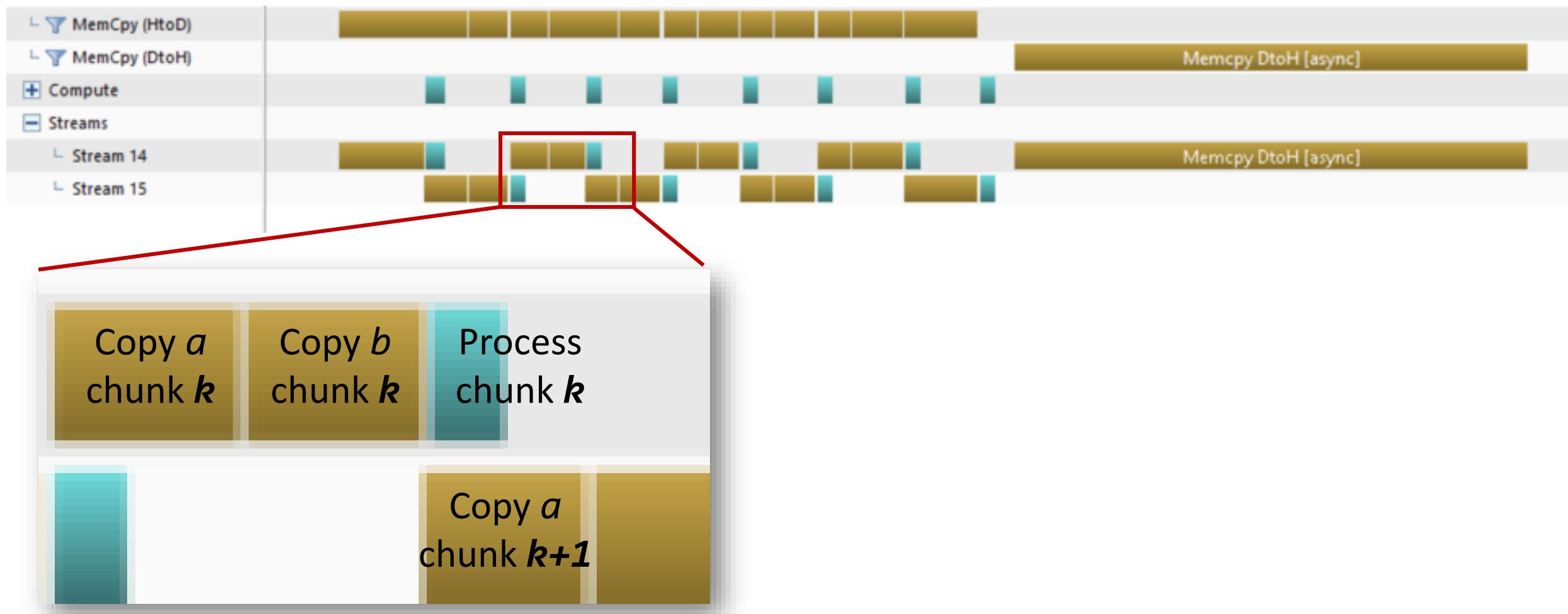


Memory in CUDA

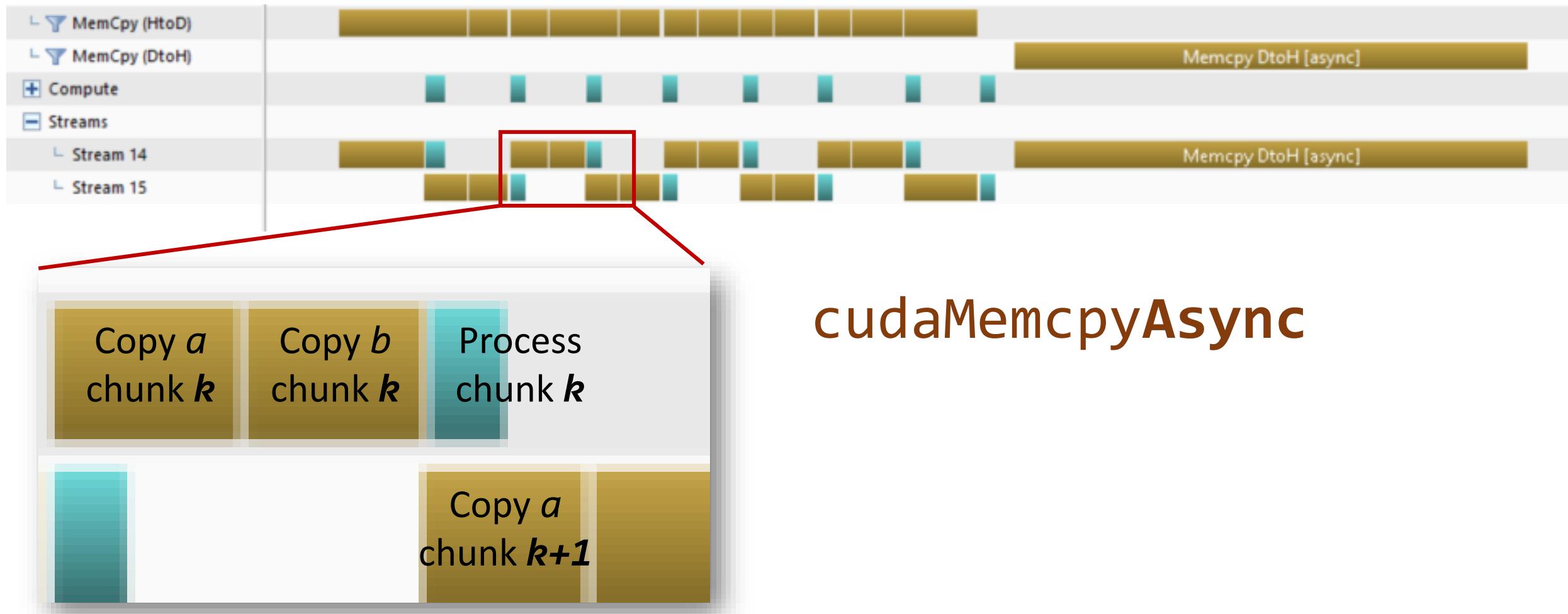
# Synchronous Copy



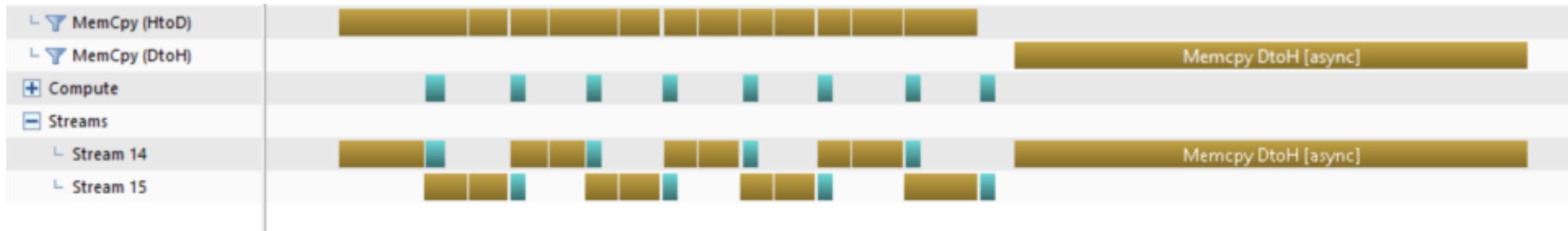
# Asynchronous Copy



# Asynchronous Copy



# Asynchronous Copy



Pinned Host Memory

**cudaMemcpyAsync**

**cudaHostAlloc**

**cudaHostRegister**

# Zero Copy Memory



# Zero Copy Memory



1. Pin Memory (`cudaHostRegisterMapped`)

`cudaHostAlloc`

`cudaHostRegister`

# Zero Copy Memory



1. Pin Memory (`cudaHostRegisterMapped`)

`cudaHostAlloc`

`cudaHostRegister`

2. `cudaHostGetDevicePointer`

3. Launch the kernel

# Zero Copy Memory



## 1. Pin Memory (`cudaHostRegisterMapped`)

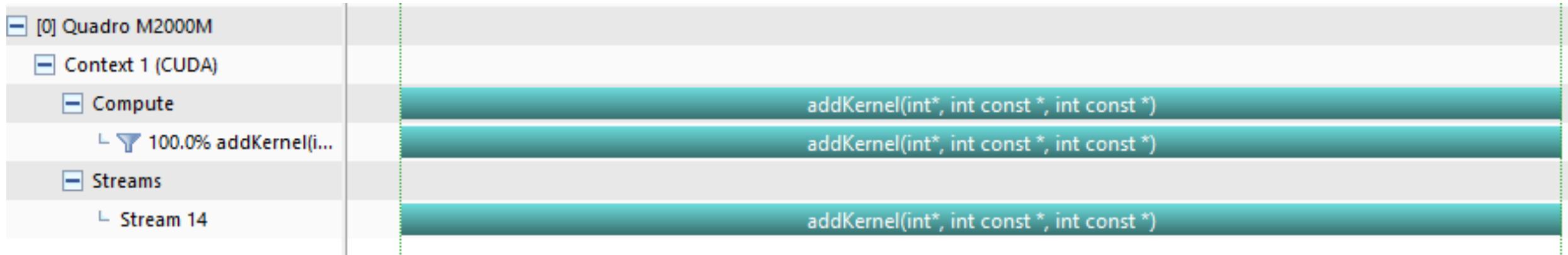
`cudaHostAlloc`

`cudaHostRegister`

## 2. ~~`cudaHostGetDevicePointer`~~

## 3. Launch the kernel

Unified Virtual Addressing (*UVA*)



**2.2 msec (Hand Made)**



**70.1 msec (Zero Copy)**



# Unified Memory

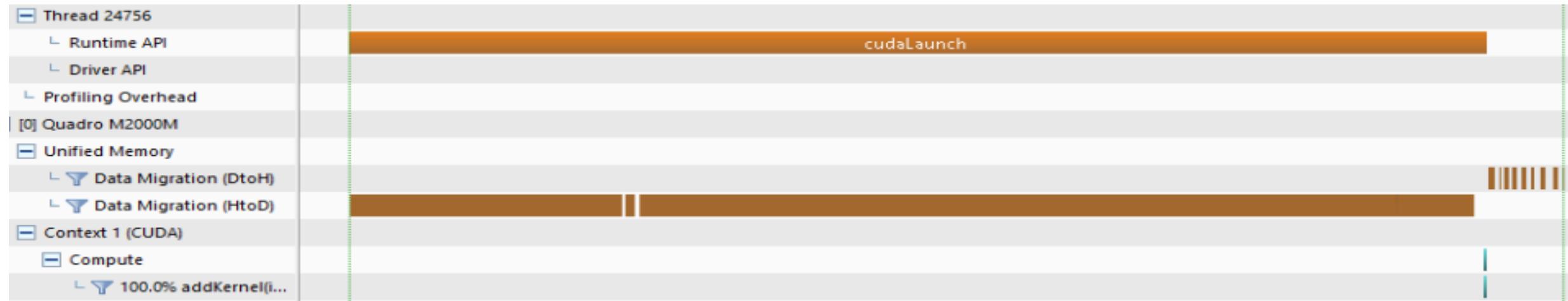
- Starting from CUDA 6
- Not to be confused with UVA!

# Unified Memory

## cudaMallocManaged

### Memory migration:

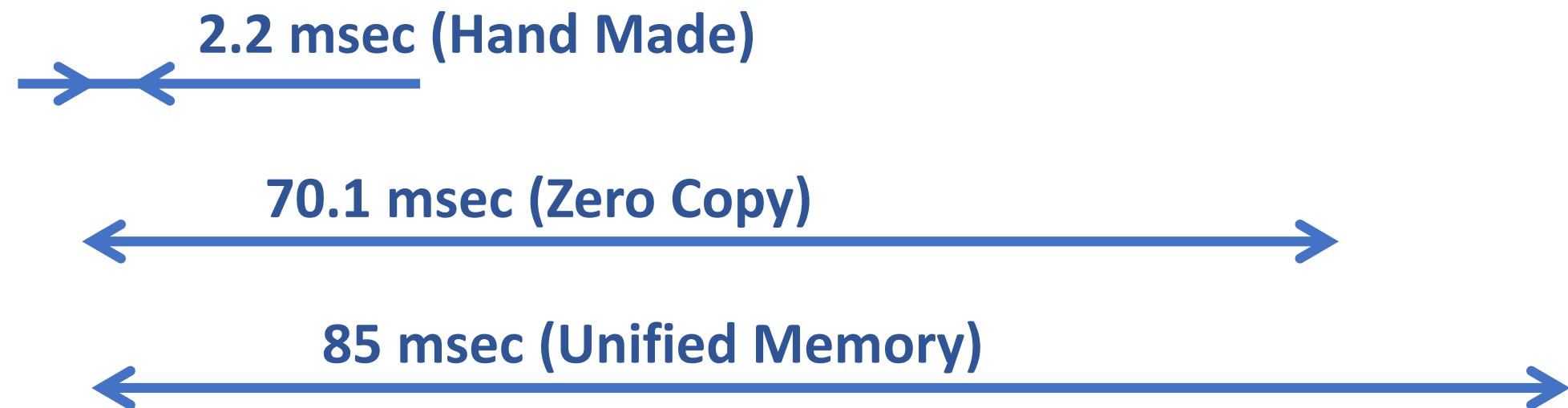
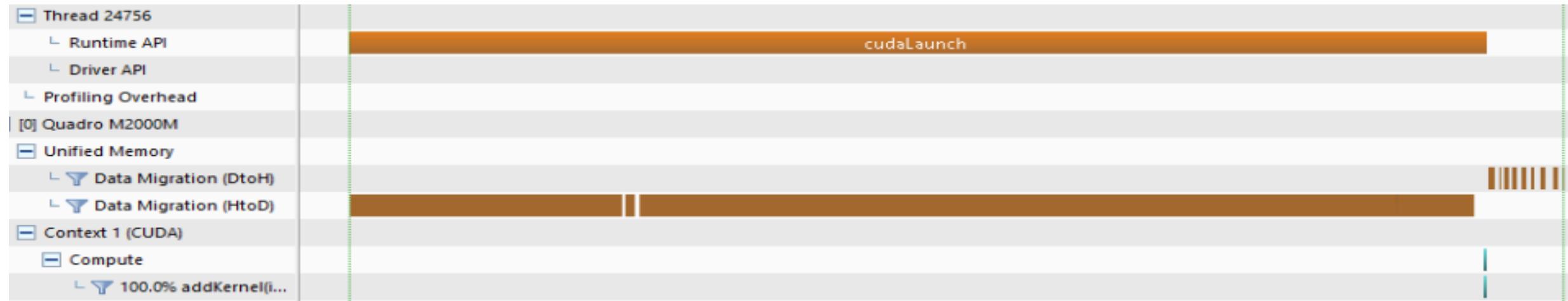
1. Data resides on the Host memory
2. Device code tries to read the memory
3. Page Fault → the driver *migrates* the memory to Device
  - Copy
4. CPU Page is swapped out



2.2 msec (Hand Made)

70.1 msec (Zero Copy)

85 msec (Unified Memory)



cudaMemPrefetchAsync, cudaMemAdvise

	Performance	Simplicity	No pinned memory	Sparse Access
Regular Copy	✓		✓	
Async Copy	✓✓			
Zero Copy	✗	✓		✓
Unified Memory	✗✓	✓✗	✓	✓



Professional  
**CUDA C**  
Programming

<static>  
**Polymorphism**

**Dynamic  
Polymorphism**

$\lambda$   
**new/delete**

# Simplest lambda

```
template<typename T>
__global__ void addKernel(DevPtr<T> c, DevPtr<const T> a, DevPtr<const T> b) {
    int idx = (blockIdx.x * blockDim.x) + threadIdx.x;
    auto op = [](__device__ auto a, __device__ auto b){ return a + b; };
    c[idx] = op(a[idx], b[idx]);
}
```

# Regular capture rules apply

```
template<typename T>
__global__ void addKernel(DevPtr<T> c, DevPtr<const T> a, DevPtr<const T> b) {
    int idx = (blockIdx.x * blockDim.x) + threadIdx.x;
    auto op = [&]{ return a[idx] + b[idx]; };
    c[idx] = op();
}
```

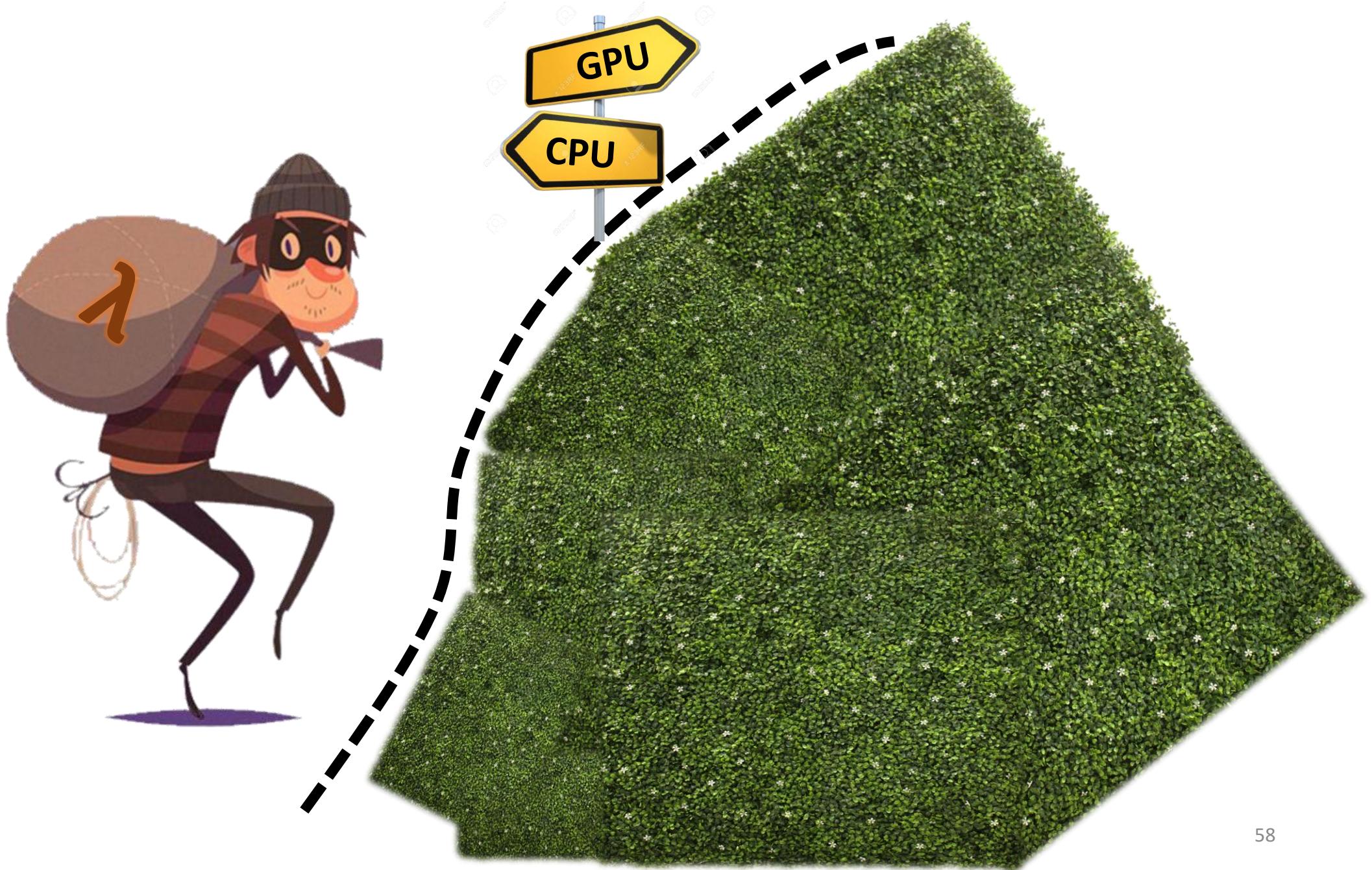
# Lambda parameters!!

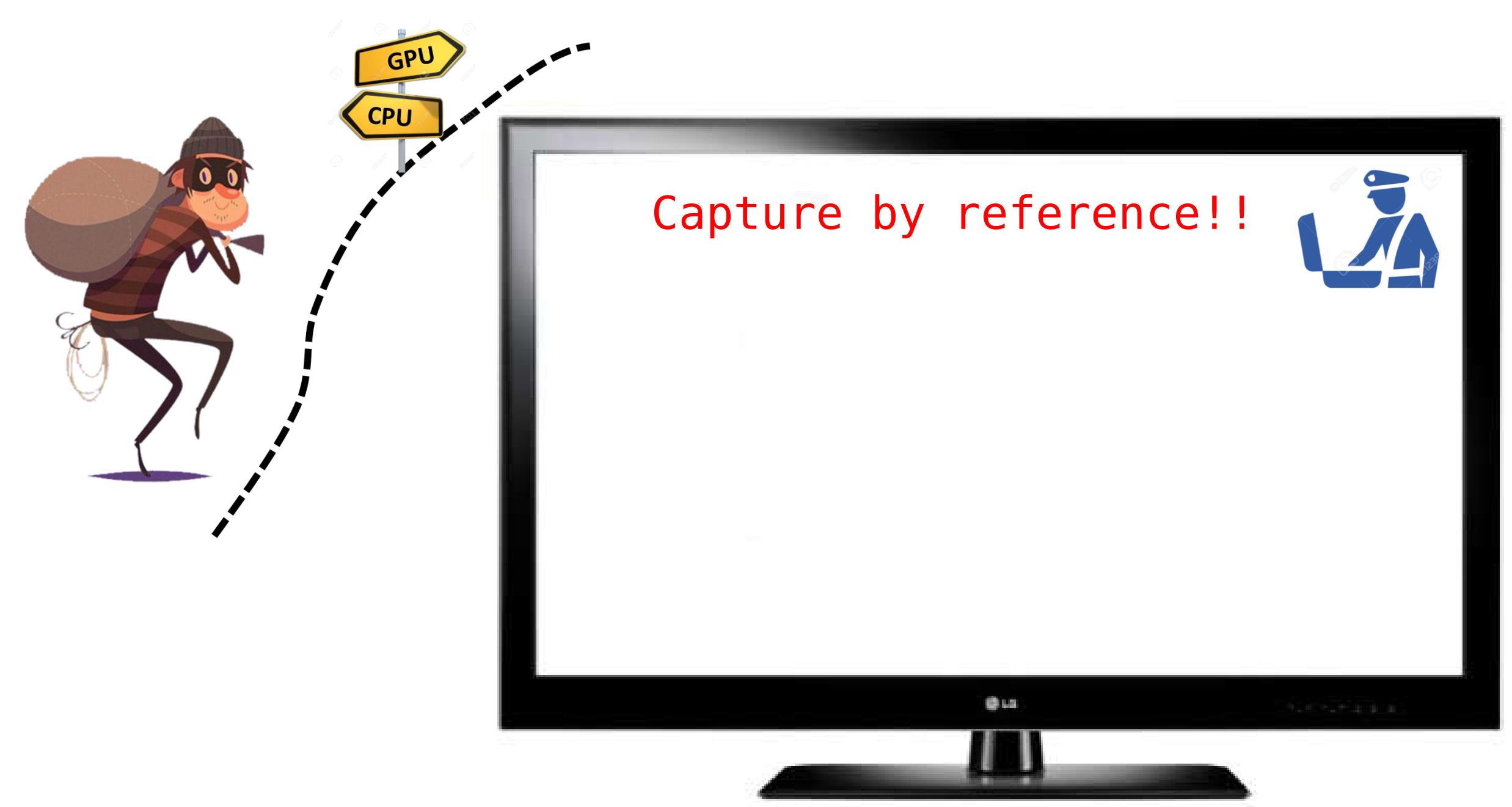
```
template<typename T, typename OP>
__global__ void addKernel(DevPtr<T> c, DevPtr<const T> a, DevPtr<const T> b,
                         OP op) {
    int idx = (blockIdx.x * blockDim.x) + threadIdx.x;
    c[idx] = op(a[idx], b[idx]);
}

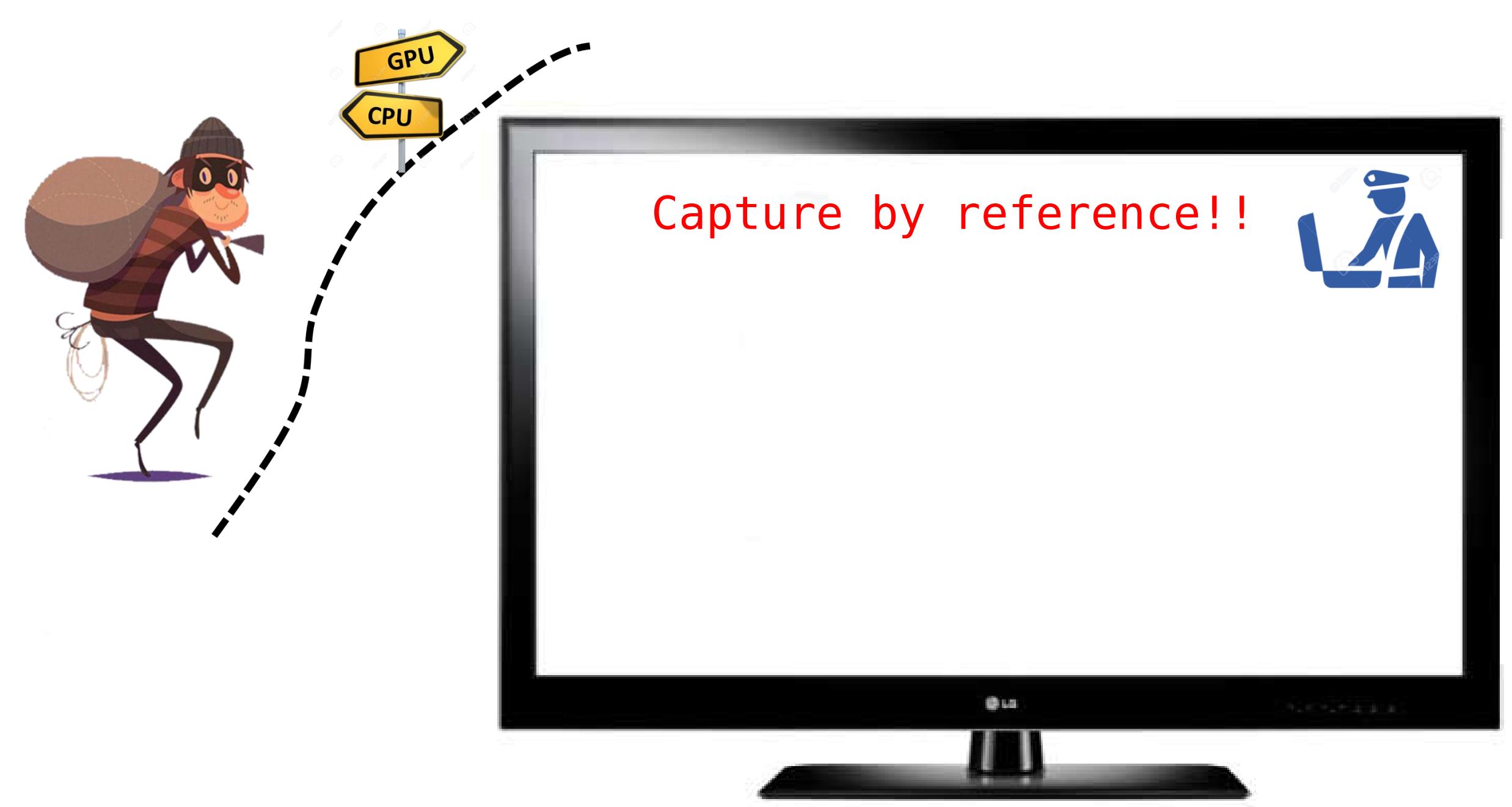
int main() {
    //...
    auto op = [] __device__ (auto a, auto b){ return a + b; };
    addKernel<<<blocks, 32>>>(cDev, aDev, bDev, op);
    //...
}
```

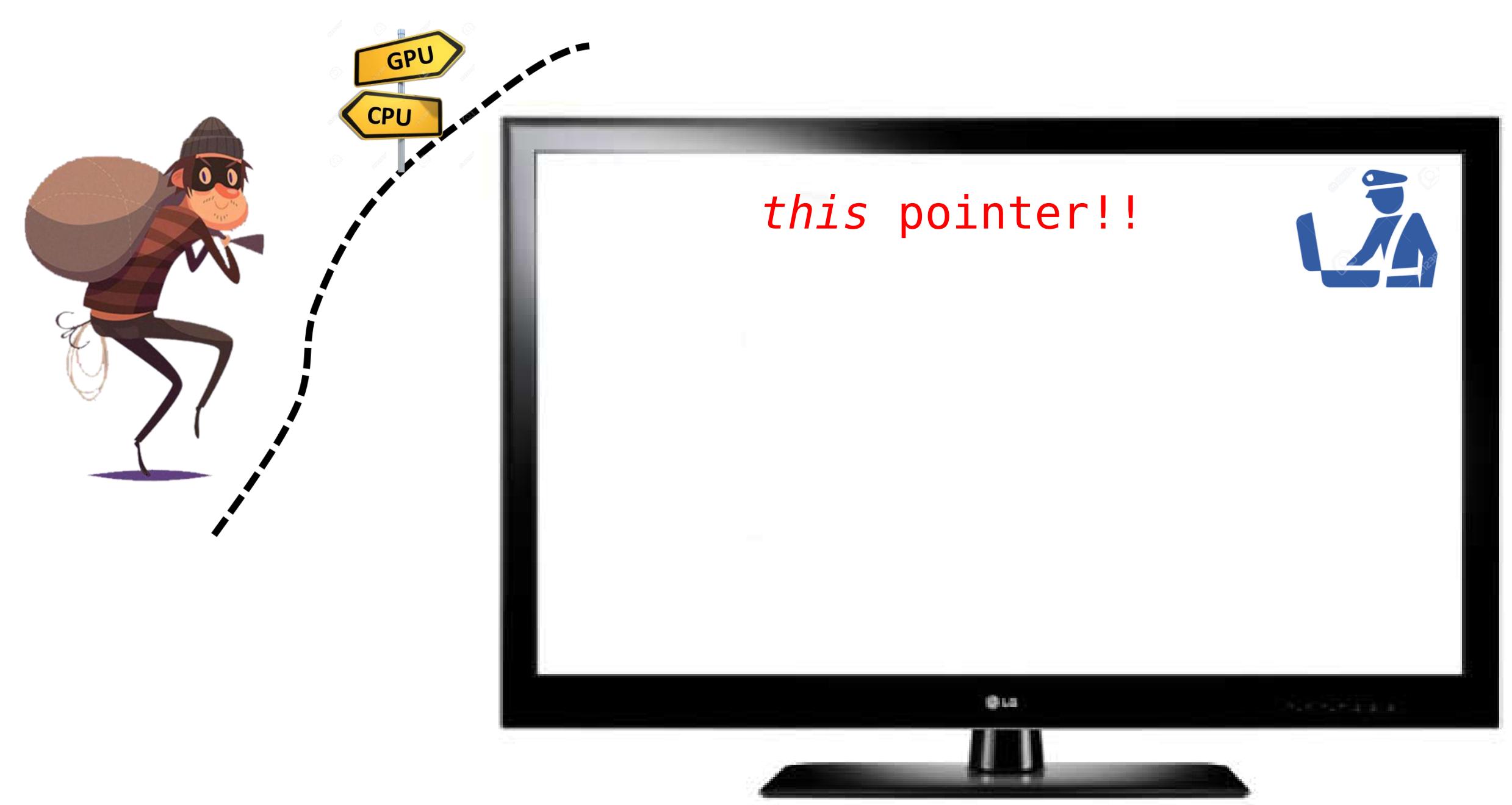
Note the *\_\_device\_\_* keyword

Requires *--expt-extended-Lambda*  
compilation flag









```
struct OP {  
    int _i;  
    explicit OP(int i) : _i(i) {}  
  
    template<typename TC, typename TAB, typename DIM>  
    void apply(TC &cDev, TAB &aDev, TAB &bDev, DIM blocks) {  
        auto op = [this] __device__ (auto a, auto b){ return a + b + _i; };  
        addKernel<<<blocks, 32>>>(cDev, aDev, bDev, op);  
    }  
};  
  
int main() {  
    //...  
    OP op{42};  
    op.apply(cDev, aDev, bDev, blocks);  
    //...  
}
```



```
Host Host
struct OP {
    int _i;
    explicit OP(int i) : _i(i) {}

    template<typename TC, typename TAB, typename DIM>
    void apply(TC &cDev, TAB &aDev, TAB &bDev, DIM blocks) {
        auto op = [*this] __device__ (auto a, auto b){ return a + b + _i; };
        addKernel<<<blocks, 32>>>(cDev, aDev, bDev, op);
    }
};

int main() {
    //...
    OP op{42};
    op.apply(cDev, aDev, bDev, blocks);
    //...
}
```



Host

```
struct OP {
    int _i;
    explicit OP(int i) : _i(i) {}

    template<typename TC, typename TAB, typename DIM>
    void apply(TC &cDev, TAB &aDev, TAB &bDev, DIM blocks) {
        auto op = [*this] __device__ (auto a, auto b){ return a + b + _i; };
        addKernel<<<blocks, 32>>>(cDev, aDev, bDev, op);
    }
};

int main() {
    //...
    OP op{42};
    op.apply(cDev, aDev, bDev, blocks);
    //...
}
```



You KNOW TOO MUCH!

```
struct OP {  
    int _i;  
    explicit OP(int i) : _i(i) {}  
  
    auto make_op() {  
        return [*this] __device__ (auto a, auto b){ return a + b + _i; };  
    }  
};  
  
int main() {  
    //...  
    OP op{42};  
    ??????  
    //...  
}
```

Host

Host

```
struct OP {  
    int _i;  
    explicit OP(int i) : _i(i) {}  
  
    auto make_op() {  
        return [*this] __device__ (auto a, auto b){ return a + b + _i; };  
    }  
};  
  
int main() {  
    //...  
    OP op{42};  
    addKernel<<<blocks, 32>>>(cDev, aDev, bDev, op.make_op());  
    //...  
}
```

### Attempt 1

error : The enclosing parent function ("make\_op") for an extended \_\_device\_\_ lambda must not have deduced return type

```
struct OP {  
    int _i;  
    explicit OP(int i) : _i(i) {}  
  
    std::function<int(int, int)> make_op() {  
        return [*this] __device__ (auto a, auto b){ return a + b + _i; };  
    }  
};  
  
int main() {  
    //...  
    OP op{42};  
    addKernel<<<blocks, 32>>>(cDev, aDev, bDev, op.make_op());  
    //...  
}
```

### Attempt 2

error : calling a \_\_host\_\_  
function("std::\_\_Func\_class<int > ::operator ()  
const") from a \_\_global\_\_ function...

```
struct OP {  
    int _i;  
    explicit OP(int i) : _i(i) {}  
  
    nvstd::function<int(int, int)> make_op() {  
        return [*this] __device__ (auto a, auto b){ return a + b + _i; };  
    }  
};  
  
int main() {  
    //...  
    OP op{42};  
    addKernel<<<blocks, 32>>>(cDev, aDev, bDev, op.make_op());  
    //...  
}
```

#include &lt;nvfunctional&gt;

Attempt 3

Compiles but fails at runtime

*... cannot be passed from host code to device code (and vice versa) at run time ...*

```
struct OP {  
    int _i;  
    explicit OP(int i) : _i(i) {}  
  
    nvstd::function<int(int, int)> __device__ __host__ make_op() {  
        return [this] (auto a, auto b){ return a + b + _i; };  
    }  
};  
  
int main() {  
    //...  
    OP op{42};  
    addKernel<<<blocks, 32>>>(cDev, aDev, bDev, op);  
    //...  
}
```

Pass the whole object to kernel, create the function using **make\_op** in the kernel

Device

Host

Professional  
**CUDA C**  
Programming

`nvstd::  
function<>`

<static>  
**Polymorphism**

**Dynamic  
Polymorphism**

`new/delete`

*But are these really Zero-Overhead  
(runtime) abstractions?*



# Godbolting CUDA



The screenshot shows the Compiler Explorer interface with three main panes:

- Left Pane (Editor):** Displays CUDA source code. A specific line of code using a lambda expression is highlighted in light blue. The code is as follows:

```
1 // Comparing function/lambda to "regular" C code,
2 // to make sure these are 0-cost abstractions
3 // See https://migocpp.wordpress.com/2018/04/02/cuda-lambdas/ for more details
4 #include <nvfunctional>
5
6 struct AddValue {
7     int _i;
8
9     AddValue(int i) : _i(i) {}
10
11    nvstd::function<int(int, int)> __device__ make_op() {
12        return [*this](auto a, auto b){ return a + b + _i; };
13    }
14};
15
16 template<typename OP>
17 __global__ void applyKernelOp(int* c, int* a, int* b, OP op)
18{
19    auto idx = threadIdx.x;
20    c[idx] = op.make_op()(a[idx], b[idx]);
21}
22
23 __global__ void applyKernelDirect(int* c, int* a, int* b, int val)
24{
25    auto idx = threadIdx.x;
26    c[idx] = a[idx] + b[idx] + val;
27}
28
29
30
31 void f() {
32     AddValue op(42);
33     applyKernelOp<<<1, 1>>>(nullptr, nullptr, nullptr, op);
```

- Middle Pane (NVCC 9.2):** Shows the generated assembly code. The highlighted line is:

```
.visible .entry _Z17applyKernelDirectPiS_i(
    .param .u64 _Z17applyKernelDirectPiS_i_param_0,
    .param .u64 _Z17applyKernelDirectPiS_i_param_1,
    .param .u64 _Z17applyKernelDirectPiS_i_param_2,
    .param .u32 _Z17applyKernelDirectPiS_i_param_3
)
```

- Bottom Pane (NVCC 9.2):** Shows the assembly code again, with the same highlighted line.

A red arrow points from the highlighted lambda expression in the editor to the corresponding assembly code in the middle pane. A large orange callout at the bottom right contains the URL: <https://godbolt.org/g/mztqWk>.

# cuobjdump

```
C:\cuda\Release> cuobjdump lambda.cu.obj -sass

...
Function : _Z17applyKernelDirectN7cudacpp12DeviceVectorIiEES1_S1_i
.headerflags      @"EF_CUDA_SM30 EF_CUDA_PTX_SM(EF_CUDA_SM30)“

/*0008*/    MOV R1, c[0x0][0x44];          /* 0x2800400110005de4 */
/*0010*/    S2R R0, SR_TID.X;            /* 0x2c0000084001c04 */
/*0018*/    MOV32I R7, 0x4;              /* 0x18000001001dde2 */
/*0020*/    ISCADD R2.CC, R0, c[0x0][0x150], 0x2; /* 0x4001400540009c43 */
...
```

Professional  
**CUDA C**  
Programming

**TAX  
FREE**  
nvstd::  
function<>

**TAX  
FREE**  
 $\lambda$

<static>  
**TAX  
FREE**  
Polymorphism

Dynamic  
Polymorphism

**new/delete**

Professional  
**CUDA C**  
Programming

nvstd::  
**function<>**



auto  
**constexpr**  
**for(a: A)**  
**A&&/std::move**

<static>  
**Polymorphism**

**Dynamic**  
**Polymorphism**

**new/delete**

# Professional CUDA C Programming



nvstd::  
**function<>**



auto  
**constexpr**  
**for(a: A)**  
**A&&/std::move**

<static>  
**Polymorphism**

**Dynamic**  
**Polymorphism**

**new/delete**

$\lambda$

# Professional CUDA C Programming



nvstd::  
function<>



auto  
constexpr  
for(a: A)  
A&&/std::move

<static>  
Polymorphism

Dynamic  
Polymorphism

new/delete

#pragma  
unroll

λ

```
template<typename T>
__device__ void myKernel(T *in, size_t length) {
    for (auto i = 0; i < length; ++i)
        // ...
}
```

```
myKernel    <<<...>>>(in, length);
```

```
template<typename T>
__device__ void myKernel(T *in, size_t length) {
    #pragma unroll ???
    for (auto i = 0; i < length; ++i)
        // ...
}
```

```
myKernel    <<<...>>>(in, length);
```

```
template<typename T>
__device__ void myKernel(T *in, size_t length) {
    #pragma unroll 10
    for (auto i = 0; i < length; ++i)
        // ...
}
```

```
myKernel    <<<...>>>(in, length);
```

Device

```
constexpr __host__ __device__ int mymax(int x, int y) {return ...}
```

```
template<int unrollFactor, typename T>
__device__ void myKernel(T *in, size_t length) {
    #pragma unroll mymax(unrollFactor, 32)
    for (auto i = 0; i < length; ++i)
        // ...
}
```

Host

```
myKernel<64><<<...>>>(in, length);
```

# CUDA Runtime Compilation

- Runtime tuning of compilation flags (*architecture* etc.)
- **Runtime selection of template parameters**

```
template<int LAYERS, typename T>
__global__ void process(T *data) {
    #pragma unroll LAYERS
    //...
}
```

```
void main() {
    int layers = /*...*/
    //...
    process<????><<<...>>>(data);
}
```

```
template<int LAYERS, typename T>
__global__ void process(T *data) { /*...*/}
```

```
void doProcess(int layers, int* data) {
    if (layers == 1) process<1><<<...>>>(data);
    if (layers == 2) process<2><<<...>>>(data);
    if //...
}

void main() {
    doProcess(layers, data);
}
```

All the template  
instantiations are being  
compiled

# Using Dynamic Compilation

```
// load the source file and headers into buffer  
nvrtcCreateProgram // from source
```

# Using Dynamic Compilation

```
// load the source file and headers into buffer  
nvrtcCreateProgram          // create from source  
nvrtcAddNameExpression    // register the instantiation
```

# Using Dynamic Compilation

```
// load the source file and headers into buffer  
nvrtcCreateProgram          // create from source  
nvrtcAddNameExpression      // register the instantiation  
  
nvrtcAddNameExpression(prog, "process<10, int>");
```

# Using Dynamic Compilation

```
// load the source file and headers into buffer
nvrtcCreateProgram          // create from source
nvrtcAddNameExpression      // register the instantiation

thrust::device_vector<int> v;
nvrtcAddNameExpression(prog, ???);
// class thrust::device_vector<int,
//   class thrust::device_malloc_allocator<int> >
```

# Using Dynamic Compilation

```
// load the source file and headers into buffer
nvrtcCreateProgram          // create from source
nvrtcAddNameExpression      // register the instantiation

thrust::device_vector<int> v;
std::string name;
nvrtcGetTypeName<decltype(v)>(&name);
nvrtcAddNameExpression(prog, /*build the string*/);
```

# Using Dynamic Compilation

```
// load the source file and headers into buffer  
nvrtcCreateProgram          // create from source  
nvrtcAddNameExpression    // register the instantiation  
(nvrtcGetTypeName)        // build the inst string
```

# Using Dynamic Compilation

```
// load the source file and headers into buffer  
nvrtcCreateProgram          // create from source  
nvrtcAddNameExpression    // register the instantiation  
  (nvrtcGetTypeName)        // build the inst string  
nvrtcCompileProgram        // compile to PTX format
```

# Using Dynamic Compilation

```
// load the source file and headers into buffer
nvrtcCreateProgram          // create from source
nvrtcAddNameExpression      // register the instantiation
    (nvrtcGetTypeName)        // build the inst string
nvrtcCompileProgram         // compile to PTX format
cuModuleLoadDataEx          // loads binary to GPU
```

# Using Dynamic Compilation

```
// load the source file and headers into buffer  
nvrtcCreateProgram          // create from source  
nvrtcAddNameExpression    // register the instantiation  
  (nvrtcGetTypeName)        // build the inst string  
nvrtcCompileProgram        // compile to PTX format  
cuModuleLoadDataEx         // loads binary to GPU  
cuModuleGetFunction        // finds the kernel
```

# Using Dynamic Compilation

```
// load the source file and headers into buffer  
nvrtcCreateProgram          // create from source  
nvrtcAddNameExpression    // register the instantiation  
  (nvrtcGetTypeName)        // build the inst string  
nvrtcCompileProgram        // compile to PTX format  
cuModuleLoadDataEx         // loads binary to GPU  
cuModuleGetFunction        // finds the kernel  
  (nvrtcGetLoweredName) // mangle C++ names
```

# Using Dynamic Compilation

```
// load the source file and headers into buffer  
nvrtcCreateProgram          // create from source  
nvrtcAddNameExpression    // register the instantiation  
  (nvrtcGetTypeName)        // build the inst string  
nvrtcCompileProgram        // compile to PTX format  
cuModuleLoadDataEx         // loads binary to GPU  
cuModuleGetFunction        // finds the kernel  
  (nvrtcGetLoweredName) // mangle C++ names  
cuLaunchKernel            // finally, LAUNCH!
```



- Template library, *inspired* by STL
- Data is represented as vectors: `device_vector`, `host_vector` + many algorithms provided on these vectors (`transform`, `sort`, ...)
- `device_ptr` is very similar to `DevicePtr` as we discussed

# Professional CUDA C Programming



nvstd::  
function<>  
  
auto  
constexpr  
for(a: A)  
A&&/std::move

<static>  
Polymorphism

Dynamic  
Polymorphism

new/delete



#pragma  
unroll

# Professional CUDA C Programming

<static>  
Polymorphism

Dynamic  
Polymorphism

nvstd::  
function<>

λ

new/delete



auto  
constexpr  
for(a: A)  
A&&/std::move

#pragma  
unroll



Using C++ !!



@michael\_gop



[mgopshtein/cudacpp](https://github.com/mgopshtein/cudacpp)  
(code examples)



<https://corecppil.github.io/Meetups/>

- New Compiler Features in CUDA 8  
<https://devblogs.nvidia.com/new-compiler-features-cuda-8/>
- Kokkos: C++ Programming model for HPC  
<https://github.com/kokkos/kokkos>