

# Efficient GPGPU programming

Ashot Vardanian

# Who am I?

Ashot Vardanian, 24  
First OpenGL line in ~15yo

Working on:

- High Performance Computing
- AI Research

Worked on:

- Web
- Mobile
- Desktop
- Scientific Computing

[github.com/ashvardanian](https://github.com/ashvardanian)  
[fb.com/ashvardanian](https://fb.com/ashvardanian)

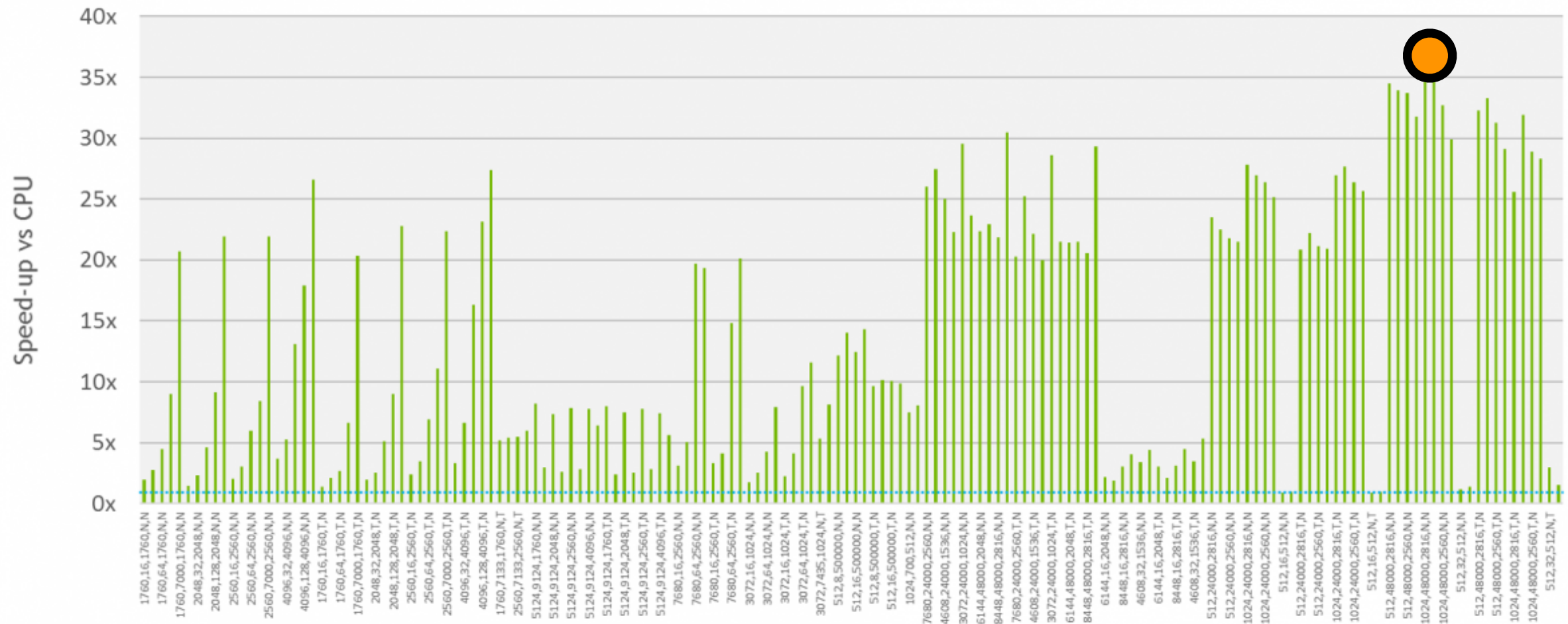


# Who is this talk for?

- You are familiar with C/C++.
- You know what a GPU is.
- You want to do number-crunching.

# Why GPUs?

...I have heard we can get a 35x performance increase...



# What we hope for?

Write code once , but

Run everywhere!

Max performance , but

Min boilerplate!

# What we hope for?

Write code once

.cpp

, but

Run everywhere!

Intel, Nvidia GPUs, AMD, Xilinx FPGA

Max performance

, but

Min boilerplate!

# What we hope for?

**Write code once**

Unified Language

**, but**

**Run everywhere!**

Modular Compilers

**Max performance**

Tune code without rewriting logic

**, but**

**Min boilerplate!**

Clean APIs

# Comparison of recipes

...we will fill this table:

	Simple	Unified	Flexible	Clean
Technology	?	?	?	?
Write code once	?	?	?	?
Run everywhere	?	?	?	?
Max performance	?	?	?	?
Minimal code size	?	?	?	?



# Plan

- 1. Popular APIs:**
  - 1. OpenGL,**
  - 2. OpenCL.**
2. Writing Low-level code
3. Existing Libraries & Tools
4. Optimal Recipes

# Plan

1. Popular APIs
2. **Writing Low-level code:**
  1. **OpenCL Language,**
  2. **CUDA Language,**
  3. **GLSL.**
3. Existing Libraries & Tools
4. Optimal Recipes

# Plan

1. Popular APIs
2. Writing Low-level code
3. **Existing Libraries & Tools:**
  1. **Linear Algebra,**
  2. **Lazy Evaluation,**
  3. **Halide,**
  4. **SyCL.**
4. Optimal Recipes

# Plan

1. Popular APIs
2. Writing Low-level code
3. Existing Libraries & Tools
4. **Optimal Recipes.**

# Popular APIs

For CPU-GPU communication

# API Support

	OpenGL
Release	1992, SGI
Intel	<b>Yes</b>
AMD	<b>Yes</b>
Nvidia	<b>Yes</b>
Apple	Deprecated
Android	<b>Yes</b>

# API Support

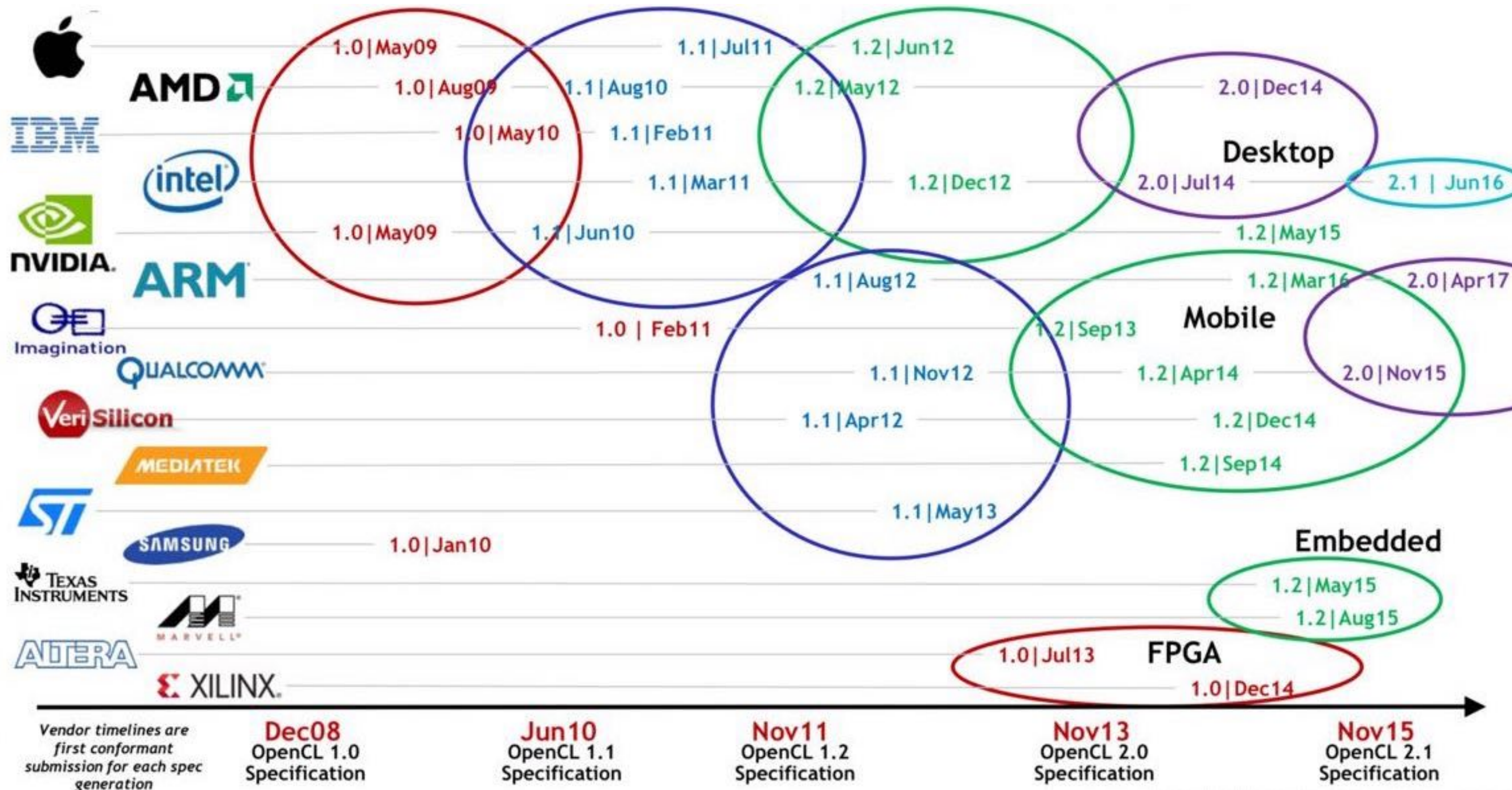
	OpenGL	CUDA
Release	1992, SGI	2007, Nvidia
Intel	<b>Yes</b>	No
AMD	<b>Yes</b>	No
Nvidia	<b>Yes</b>	<b>Yes</b>
Apple	Deprecated	No
Android	<b>Yes</b>	No

# API Support

	OpenGL	CUDA	OpenCL
Release	1992, SGI	2007, Nvidia	2009, Apple
Intel	<b>Yes</b>	No	<b>Yes</b>
AMD	<b>Yes</b>	No	<b>Yes</b>
Nvidia	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>
Apple	Deprecated	No	MacOS
Android	<b>Yes</b>	No	Depends



# OpenCL support (2015)



# API Support

	OpenGL	CUDA	OpenCL	Metal
Release	1992, SGI	2007, Nvidia	2009, Apple	2014, Apple
Intel	<b>Yes</b>	No	<b>Yes</b>	No
AMD	<b>Yes</b>	No	<b>Yes</b>	No
Nvidia	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>	No
Apple	Deprecated	No	MacOS	<b>Yes</b>
Android	<b>Yes</b>	No	Depends	No

# API Support

	OpenGL	CUDA	OpenCL	Metal	Vulkan
Release	1992, SGI	2007, Nvidia	2009, Apple	2014, Apple	2016, AMD
Intel	<b>Yes</b>	No	<b>Yes</b>	No	<b>Yes</b>
AMD	<b>Yes</b>	No	<b>Yes</b>	No	<b>Yes</b>
Nvidia	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>	No	<b>Yes</b>
Apple	Deprecated	No	MacOS	<b>Yes</b>	MoltenVK
Android	<b>Yes</b>	No	Depends	No	<b>Yes</b>

# API Comparison

	OpenGL	CUDA	OpenCL	Metal	Vulkan
Primary Purpose	Graphics	<b>Compute</b>	<b>Compute</b>	Graphics	Graphics
Base Input Language	C	<b>C++</b>	C	<b>C++</b>	<b>Any</b>
Complexity	Hard ...on Device	<b>Easy</b>	<b>Easy</b>	Average	Very Hard
Targets Flexibility	Average	Low ...only Nvidia	<b>Extreme ...FPGA</b>	Low ...only Apple	<b>High</b>
API Flexibility*	Average	<b>High**</b>	Average	Average	<b>High</b>

# API Comparison

		CUDA	OpenCL		Vulkan
Primary Purpose		Compute	Compute		Graphics
Base Input Language		C++	C		Any
Complexity		Easy	Easy		Very Hard
Targets Flexibility		Low ...only Nvidia	Extreme ...FPGA		High
API Flexibility*		High**	Average		High

# Language Syntax

CUDA vs OpenCL

# Parallelism in Language

Which keywords and features must a language have to make parallel programming easy?

**Synchronization  
Primitives**



To help threads  
understand their role

**Memory  
Qualifiers**

To limit data  
visibility

?

# Memory Types

	CUDA	OpenCL
 All Threads	Global	Global
Group of Threads	Shared	Local
Single Thread	Local, Register (faster)	Private
 Other	Constant, Texture	Constant

OpenGL	Vertex Buffer	Frame Buffer	Texture	Local
--------	---------------	--------------	---------	-------



# Actual Memory Types

...have little to do with physical capabilities of the device! At least from OpenCL perspective!

	i7-7820HQ	Titan V	Radeon Pro 560
Compute Units	8 cores	80 cores	16 cores
Sync-able Group	<1024 threads $\boxed{N^1}$	<1024 threads $\boxed{N^3}$	< 256 threads $\boxed{N^3}$
Constant Buffer	64 Kb	? Kb	64 Kb
Local Memory	32 Kb	? Kb	32 Kb

1 Mb L2

16 Kb L1  
per CU

# Actual Memory Types

...have little to do with physical capabilities of the device! At least from OpenCL perspective!

	i7-7820HQ	Titan V	Radeon Pro 560
Compute Units	8 cores	80 cores	16 cores
Sync-able Group	<1024 threads $N^1$	<1024 threads $N^3$	< 256 threads $N^3$
Constant Buffer	64 Kb	"In Volta the L1 cache, texture cache, and shared memory are backed by a combined 128 KB data cache."	64 Kb
Local Memory	32 Kb		32 Kb

1 Mb L2

16 Kb L1 per CU

Nvidia GPUs have one real "constant" buffer (64-128 Kb) and allocate rest in global memory.

AMD GPUs often have multiple "constant" buffers (64 Kb each) and allocate rest in global memory.

# Memory Qualifiers

	CUDA	OpenCL
All Threads	__device__	__global
Group of Threads	__shared__	__local
Single Thread	~	~
Other	__constant__	__constant

CPU

```
void sum_2_vecs(float const * xA,
               float const * xB,
               float * y,
               int const xLen);
```

GPU version

```
kernel
void sum_2_vecs(global float const * xA,
               global float const * xB,
               global float * y);
```

# Terminology

CUDA		OpenCL
Stream Multiprocessor	Core on CPU	Compute Unit
Thread	Thread on CPU	Work-Item
Block		Work-Group
__global__ function		__kernel function
__device__ function		~

# Kernels Indexing

	CUDA	OpenCL
	gridDim	get_num_groups()
	blockDim	get_local_size()
	blockIdx	get_group_id()
	threadIdx	get_local_id()
Ugly	$\text{blockIdx} * \text{blockDim} + \text{threadIdx}$	get_global_id()
Ugly	$\text{gridDim} * \text{blockDim}$	get_global_size()

# Kernels Synchronization

	CUDA	OpenCL
Group	<code>__syncthreads()</code>	<code>barrier(...)</code>
All	<code>__threadfence()</code>	~
Group Memory	<code>__threadfence_block()</code>	<code>mem_fence(...)</code>
	~	<code>read_mem_fence(...)</code> ?
	~	<code>write_mem_fence(...)</code> ?

# Parallelism in Language

Which keywords and features must a language have to make parallel programming easy?

Synchronization  
Primitives

`gridDim`

`get_group_id(0)`

`__threadfence_block()`

Memory  
Qualifiers

`in`

`__global`

`__shared__`

?

# Code Examples

Why would you want to write low-level kernels?



# Data-Parallel Tasks

...brute-force scaling of simple  
non-concurrent problems

inputs:																
operator:	sin				exp				cos				log			
outputs:																

# Data-Parallel Tasks

...brute-force scaling of simple  
non-concurrent problems

inputs:																
inputs:																
operator:	+   -   x   ÷				pow				fmod				atan2			
outputs:																

# Vector Sum: C

```
void sum_2_vectors(float const * xA,  
                  float const * xB,  
                  float * y,  
                  int const xLen) {  
    for (int i = 0; i < xLen; i ++)  
        y[i] = xA[i] + xB[i];  
}
```

# Vector Sum: OpenCL

```
kernel void sum_2_vectors(global float const * xA,  
                           global float const * xB,  
                           global float * y) {  
    int i = get_global_id(0);  
    y[i] = xA[i] + xB[i];  
}
```

# Vector Sum: GLSL

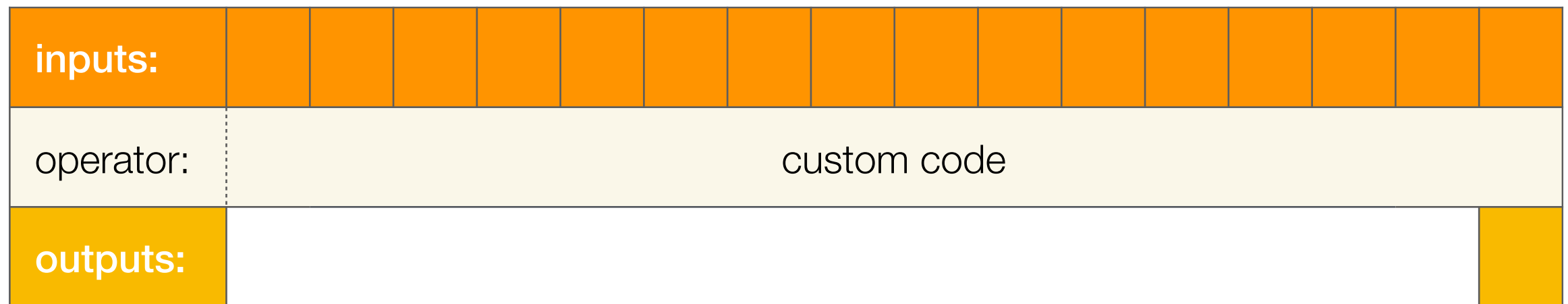
```
#version 450
```

```
layout(binding = 0) in buffer lay0 { float xA[]; };  
layout(binding = 1) in buffer lay1 { float xB[]; };  
layout(binding = 2) out buffer lay2 { float y[]; };
```

```
void main() {  
    uint const i = gl_GlobalInvocationID.x;  
    y[i] = xA[i] + xB[i];  
}
```

# Concurrent Tasks

...synchronization nightmare  
and benchmarks heaven!

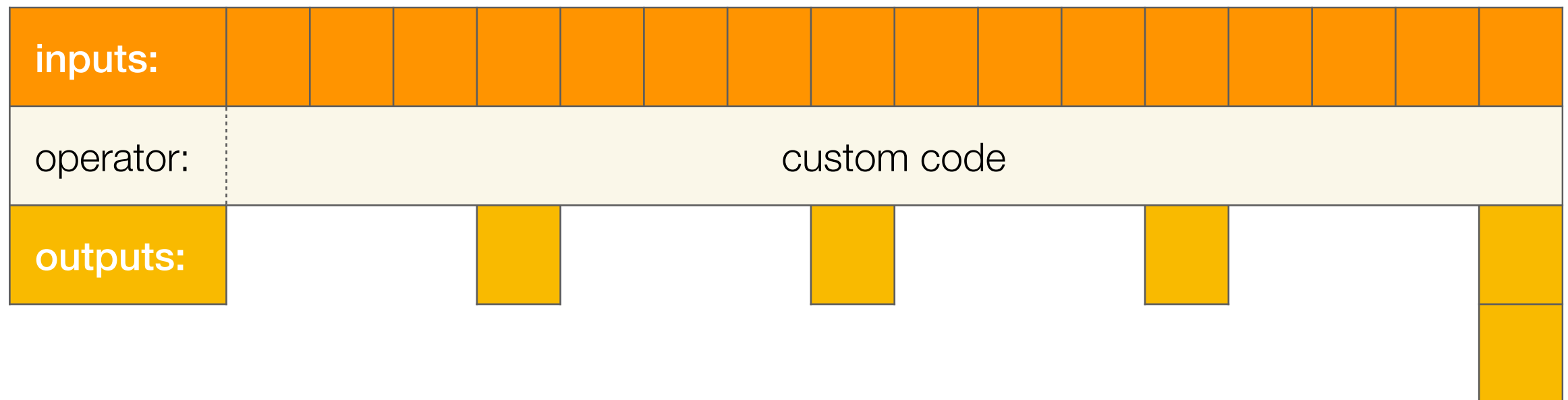


# Reduction: C

```
void reduce(float const * x,  
            float * y,  
            int const xLen) {  
    *y = 0;  
    for (int i = 0; i < xLen; i ++)  
        *y += x[i];  
}
```

# Concurrent Tasks



...force us to inject memory synchronization barriers and loops, that compiler won't unroll!





# Reduction: OpenCL (1)

```
__kernel
void reduce_simple(__global float const * xArr, __global float * yArr,
                   int const xLen, __local float * mBuffer) {
    int const lIdxGlobal = get_global_id(0);
    int const lIdxInBlock = get_local_id(0);
    mBuffer[lIdxInBlock] = (lIdxGlobal < xLen) ? xArr[lIdxGlobal] : 0;

     barrier(CLK_LOCAL_MEM_FENCE);
    int lBlockSize = get_local_size(0);
    int lBlockSizeHalf = lBlockSize / 2;
    while (lBlockSizeHalf > 0) {
        if (lIdxInBlock < lBlockSizeHalf) {
            mBuffer[lIdxInBlock] += mBuffer[lIdxInBlock + lBlockSizeHalf];
            if ((lBlockSizeHalf * 2) < lBlockSize) {
                if (lIdxInBlock == 0)
                    mBuffer[lIdxInBlock] += mBuffer[lIdxInBlock + (lBlockSize - 1)];
            }
        }
         barrier(CLK_LOCAL_MEM_FENCE);
        lBlockSize = lBlockSizeHalf;
        lBlockSizeHalf = lBlockSize / 2;
    }

    if (lIdxInBlock == 0) yArr[get_group_id(0)] = mBuffer[0];
}
```

# Reduction: OpenCL (2)

```
__kernel
void reduce_unrolled(__global float const * xArr, __global float * yArr,
                    int const xLen, __local float * mBuffer) {
    int const lIdxInBlock = get_local_id(0);
    int const lIdxGlobal = get_group_id(0) * (get_local_size(0) * 2) + get_local_id(0);
    int const lBlockSize = get_local_size(0);
    mBuffer[lIdxInBlock] = (lIdxGlobal < xLen) ? xArr[lIdxGlobal] : 0;

    if (lIdxGlobal + get_local_size(0) < xLen)
        mBuffer[lIdxInBlock] += xArr[lIdxGlobal + get_local_size(0)];
    barrier(CLK_LOCAL_MEM_FENCE);

    #pragma unroll 1
    for (int lTemp = get_local_size(0) / 2; lTemp > 32; lTemp >= 1) {
        if (lIdxInBlock < lTemp)
            mBuffer[lIdxInBlock] += mBuffer[lIdxInBlock + lTemp];
        barrier(CLK_LOCAL_MEM_FENCE);
    }

    if (lIdxInBlock < 32) {
        if (lBlockSize >= 64) { mBuffer[lIdxInBlock] += mBuffer[lIdxInBlock + 32]; }
        if (lBlockSize >= 32) { mBuffer[lIdxInBlock] += mBuffer[lIdxInBlock + 16]; }
        if (lBlockSize >= 16) { mBuffer[lIdxInBlock] += mBuffer[lIdxInBlock + 8]; }
        if (lBlockSize >= 8) { mBuffer[lIdxInBlock] += mBuffer[lIdxInBlock + 4]; }
        if (lBlockSize >= 4) { mBuffer[lIdxInBlock] += mBuffer[lIdxInBlock + 2]; }
        if (lBlockSize >= 2) { mBuffer[lIdxInBlock] += mBuffer[lIdxInBlock + 1]; }
    }

    if (lIdxInBlock == 0) yArr[get_group_id(0)] = mBuffer[0];
}
```

# Existing Libs & Tools

The complexity of Choice

# Linear Algebra

	Intel MKL	cuBLAS	CLBlast
Types	Basic	<b>Basic, FP16, INT8</b>	Basic, FP16
Performance	+	<b>+++</b>	++
APIs	BLAS, LAPACK...	BLAS +	BLAS
BLAS Levels	Vector-Vector	Matrix-Vector	Matrix-Matrix
LAPACK	Least Squares	Eigenvalues	Factorization

**Optimized kernels are chained into slow pipelines!**

# Linear Algebra

	Intel MKL	cuBLAS	CLBlast
Types	Basic	<b>Basic, FP16, INT8</b>	Basic, FP16
Performance	+	<b>+++</b>	++
APIs	BLAS, LAPACK...	BLAS +	BLAS
BLAS Levels	Vector-Vector	Matrix-Vector	Matrix-Matrix
LAPACK	Least Squares	Eigenvalues	Factorization

Vectors Sum?

Array Sum Reduction?

152 Ops!

Optimized kernels are chained into slow pipelines!

# Linear Algebra

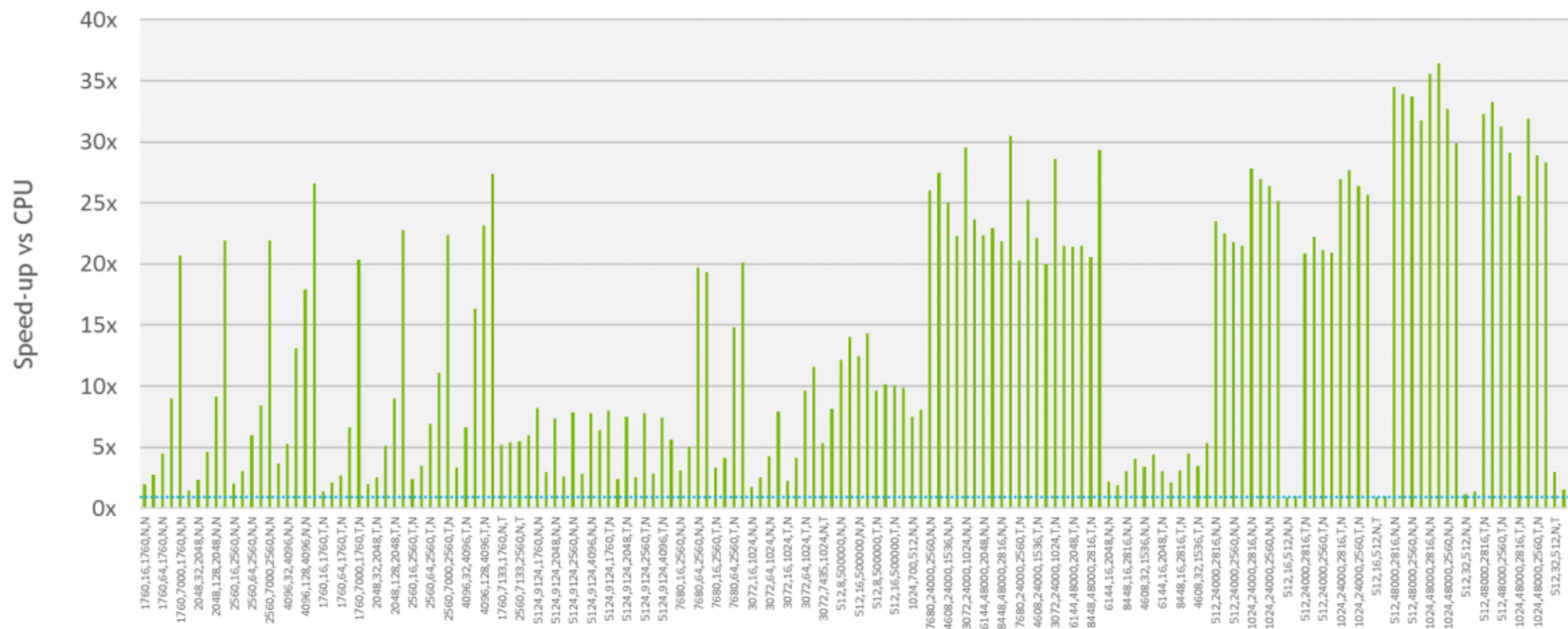
	Intel MKL	cuBLAS	CLBlast
Types	Basic	<b>Basic, FP16, INT8</b>	Basic, FP16
Performance	+	<b>+++</b>	++
APIs	BLAS, LAPACK...	BLAS +	BLAS
BLAS Levels	Vector-Vector	Matrix-Vector	Matrix-Matrix
LAPACK	Least Squares	Eigenvalues	Factorization

Vectors Sum?  
SAXPY with a=1

Array Sum Reduction?  
SDOT with unit vector

152  
Ops!

Optimized kernels are chained into slow pipelines!



	E5-2690v4	Gold 6262V	V100
Float Performance	+	++	+++
Cores	14	24	14
Year	2016	2019	2017
Price	2,000-2,500 USD	<b>3,000 USD</b>	8,000 USD

# Lazy Evaluation Graph

Lazy	Eigen	ArrayFire	Boost. Compute	Thrust	VexCL
Stars	10k	2.8k	1K	2.5k	565
Type-Safe	<b>Yes</b>	No	Yes	Yes	Yes
Backends	OpenMP, CUDA?	<b>OpenCL, CUDA, etc.</b>	OpenCL	CUDA, OpenMP	OpenCL, CUDA, OpenMP

**Very different functionality and inconsistent APIs.**  
**Potential Licensing issues.**



# Data-Parallel Tasks

...again, but now with higher level heterogeneous computing tools!

inputs:																
operator:	sin				exp				cos				log			
outputs:																

# Cost of Memory Access

...is much higher, than cost of compute, so we need kernel fusion!

	Power	
ALU	1 pJ	
Load from SRAM	3 pJ	
Move 10 mm on-chip	30 pJ	
Send off-chip	500 pJ	
Send to DRAM	1 nJ	1,000x more
Send over LTE	10 $\mu$ J	10,000,000x more

# Parallelism in Language

...we want to separate the inner part of the "for" loop and the enumeration order

Synchronization  
Primitives

To help threads  
understand their role

Memory  
Qualifiers

To limit data  
visibility

Order  
Descriptors

To simplify loops  
optimization

# How Halide works?

...by separating the inner loop logic!

```
func(i) = lA(i) + lB(i);
```

# How Halide works?

...and by making loops implicit!

```
void sum_2_vectors(float const * xA,  
                  float const * xB,  
                  float * y,  
                  int const xLen) {  
    Halide::Buffer<bFlt32> lA { const_cast<float *>(xA), xLen, "xA" };  
    Halide::Buffer<bFlt32> lB { const_cast<float *>(xB), xLen, "xB" };  
    Halide::Var i { "i" };  
    Halide::Func func;  
  
    func(i) = lA(i) + lB(i);  
  
    Halide::Buffer<bFlt32> lOut = func.realize(xLen);  
    std::copy_n(lOut.data(), xLen, y);  
}
```

Function body

The "for" loop

# How Halide works?

```
void sum_2_vectors(float const * xA,  
                  float const * xB,  
                  float * y,  
                  int const xLen) {  
    Halide::Buffer<bFlt32> lA { const_cast<float *>(xA), xLen, "xA" };  
    Halide::Buffer<bFlt32> lB { const_cast<float *>(xB), xLen, "xB" };  
    Halide::Var i { "i" };  
    Halide::Func func;  
  
    func(i) = lA(i) + lB(i);  
  
    Halide::Buffer<bFlt32> lOut = func.parallel(i).realize(xLen);  
    std::copy_n(lOut.data(), xLen, y);  
}
```



Parallel "for" loop

# How Halide works?

```
void sum_2_vectors(float const * xA,  
                  float const * xB,  
                  float * y,  
                  int const xLen) {  
    Halide::Buffer<bFlt32> lA { const_cast<float *>(xA), xLen, "xA" };  
    Halide::Buffer<bFlt32> lB { const_cast<float *>(xB), xLen, "xB" };  
    Halide::Var i { "i" };  
    Halide::Func func;  
  
    func(i) = lA(i) + lB(i);  
  
    Halide::Buffer<bFlt32> lOut = func.vectorize(i, 8).realize(xLen);  
    std::copy_n(lOut.data(), xLen, y);  
}
```



Vectorized "for" loop with "float8"

# How Halide works?

```
void sum_2_vectors(float const * xA,  
                  float const * xB,  
                  float * y,  
                  int const xLen) {  
    Halide::Buffer<bFlt32> lA { const_cast<float *>(xA), xLen, "xA" };  
    Halide::Buffer<bFlt32> lB { const_cast<float *>(xB), xLen, "xB" };  
    Halide::Var i { "i" }, j { "j" }, k { "k" };  
    Halide::Func func;  
  
    func(i) = lA(i) + lB(i);  
  
    func.vectorize(i, j, k, 8);  
    Halide::Buffer<bFlt32> lOut = func.parallel(j).unroll(k).realize(xLen);  
    std::copy_n(lOut.data(), xLen, y);  
}
```

The diagram illustrates the transformation of a 1D loop into a 2D loop and the unrolling of the inner loop. A box labeled "Transforming a 1 dimensional 'for'-loop into 2D loop" has two arrows pointing to the `parallel(j)` and `unroll(k)` methods in the code. Another box labeled "Unroll the inner loop!" has an arrow pointing to the `unroll(k)` method.



# Blur Filter: C++

...the baseline for comparison!

```
void box_filter_3x3(const Image &in, Image &blury) {  
    Image blurx(in.width(), in.height()); // allocate blurx array  
  
    for (int y = 0; y < in.height(); y++)  
        for (int x = 0; x < in.width(); x++)  
            blurx(x, y) = (in(x-1, y) + in(x, y) + in(x+1, y))/3;  
  
    for (int y = 0; y < in.height(); y++)  
        for (int x = 0; x < in.width(); x++)  
            blury(x, y) = (blurx(x, y-1) + blurx(x, y) + blurx(x, y+1))/3;  
}
```

Slow

Fast

# Blur Filter: Halide

## Halide

0.9 ms/megapixel

```
Func box_filter_3x3(Func in) {
    Func blurx, blurry;
    Var x, y, xi, yi;

    // The algorithm - no storage, order
    blurx(x, y) = (in(x-1, y) + in(x, y) + in(x+1, y))/3;
    blurry(x, y) = (blurx(x, y-1) + blurx(x, y) + blurx(x, y+1))/3;

    // The schedule - defines order, locality; implies storage
    blurry.tile(x, y, xi, yi, 256, 32)
        .vectorize(xi, 8).parallel(y);
    blurx.compute_at(blurry, x).store_at(blurry, x).vectorize(x, 8);

    return blurry;
}
```

Sugar: tiling!

## C++

0.9 ms/megapixel

With platform-specific SIMD!

```
void box_filter_3x3(const Image &in, Image &blurry) {
    __m128i one_third = _mm_set1_epi16(21846);
    #pragma omp parallel for
    for (int yTile = 0; yTile < in.height(); yTile += 32) {
        __m128i a, b, c, sum, avg;
        __m128i blurx[(256/8)*(32+2)]; // allocate tile blurx array
        for (int xTile = 0; xTile < in.width(); xTile += 256) {
            __m128i *blurxPtr = blurx;
            for (int y = -1; y < 32+1; y++) {
                const uint16_t *inPtr = &(in[yTile+y][xTile]);
                for (int x = 0; x < 256; x += 8) {
                    a = _mm_loadu_si128((__m128i*)(inPtr-1));
                    b = _mm_loadu_si128((__m128i*)(inPtr+1));
                    c = _mm_load_si128((__m128i*)(inPtr));
                    sum = _mm_add_epi16(_mm_add_epi16(a, b), c);
                    avg = _mm_mulhi_epi16(sum, one_third);
                    _mm_store_si128(blurxPtr++, avg);
                    inPtr += 8;
                }
            }
            blurxPtr = blurx;
            for (int y = 0; y < 32; y++) {
                __m128i *outPtr = (__m128i *)(&(blurry[yTile+y][xTile]));
                for (int x = 0; x < 256; x += 8) {
                    a = _mm_load_si128(blurxPtr+(2*256)/8);
                    b = _mm_load_si128(blurxPtr+256/8);
                    c = _mm_load_si128(blurxPtr++);
                    sum = _mm_add_epi16(_mm_add_epi16(a, b), c);
                    avg = _mm_mulhi_epi16(sum, one_third);
                    _mm_store_si128(outPtr++, avg);
                }
            }
        }
    }
}
```

# Laplacian Filter

...real life example!

	Reference C++
LOC	300
Time	?
Performance	1x

# Laplacian Filter

...real life example!

	Reference C++	Adobe CPU
LOC	300	1500
Time	?	3 months
Performance	1x	10x

# Laplacian Filter

...real life example!

	Reference C++	Adobe CPU	Halide CPU
LOC	300	1500	<b>60</b>
Time	?	3 months	<b>1 day</b>
Performance	1x	10x	<b>20x</b>

# Laplacian Filter

...real life example!

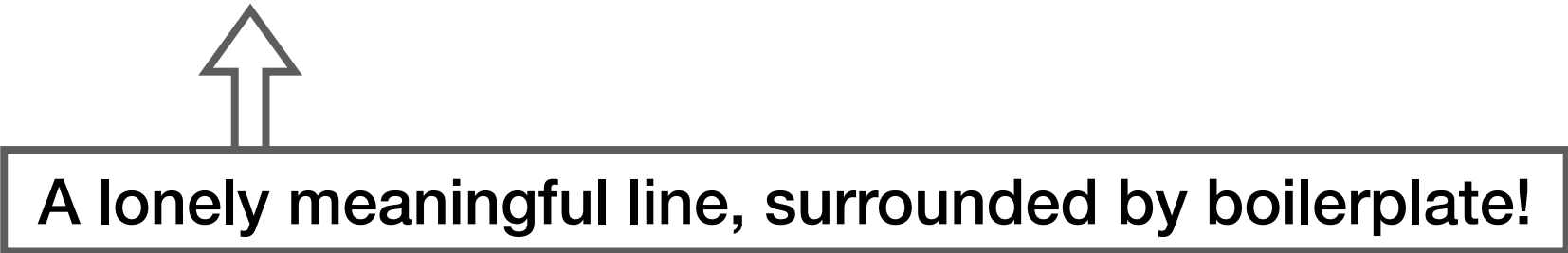
	Reference C++	Adobe CPU	Halide CPU	Halide GPU
LOC	300	1,500	<b>60</b>	
Time	?	3 months	<b>1 day</b>	
Performance	1x	10x	<b>20x</b>	<b>70x</b>

# Vector Sum: SyCL Today

```
void sum_2_vectors(float const * xA,
                  float const * xB,
                  float * y,
                  int const xLen) {

    cl::sycl::queue q;
    cl::sycl::buffer<float, 1> lA { xA, xLen };
    cl::sycl::buffer<float, 1> lB { xB, xLen };
    cl::sycl::buffer<float, 1> lOut { y, xLen };

    q.submit( [&](cl::sycl::handler & h) {
        auto hA = lA.get_access<cl::sycl::access::mode::read>(h);
        auto hB = lB.get_access<cl::sycl::access::mode::read>(h);
        auto hOut = lOut.get_access<cl::sycl::access::mode::write>(h);
        h.parallel_for<class kernel_name>(xLen, [=] (cl::sycl::id<1> i) {
            hOut[i] = hA[i] + hB[i];
        });
    });
    q.wait();
}
```



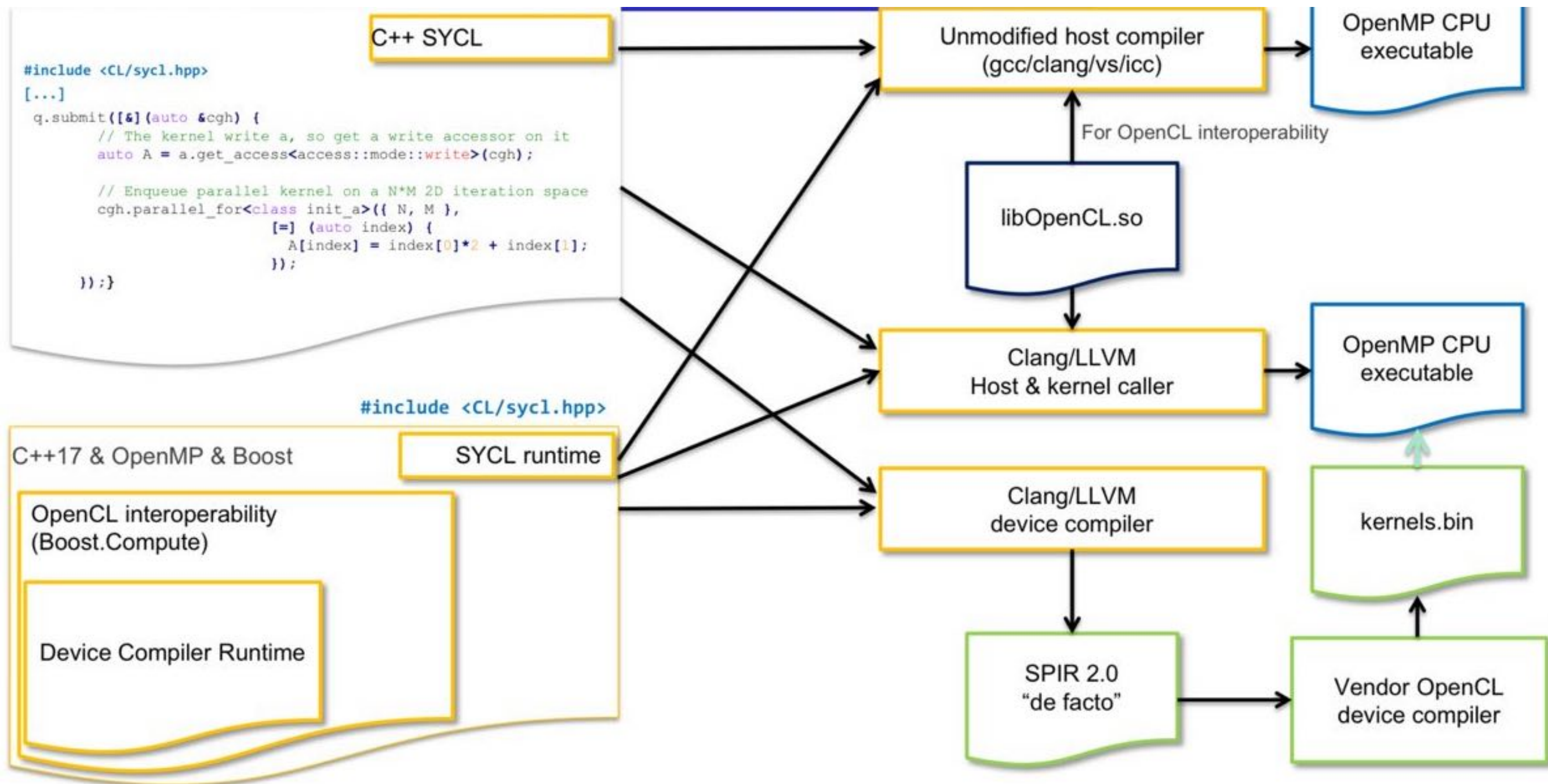
A lonely meaningful line, surrounded by boilerplate!

# Vector Sum: SyCL STL

```
void sum_2_vectors(float const * xA,  
                  float const * xB,  
                  float * y,  
                  int const xLen) {  
  
    std::transform(cl::sycl::uniform_policy,  
                  std::span(xA, xLen), std::span(xB, xLen),  
                  std::span(y, xLen),  
                  std::plus<float> { });  
}
```

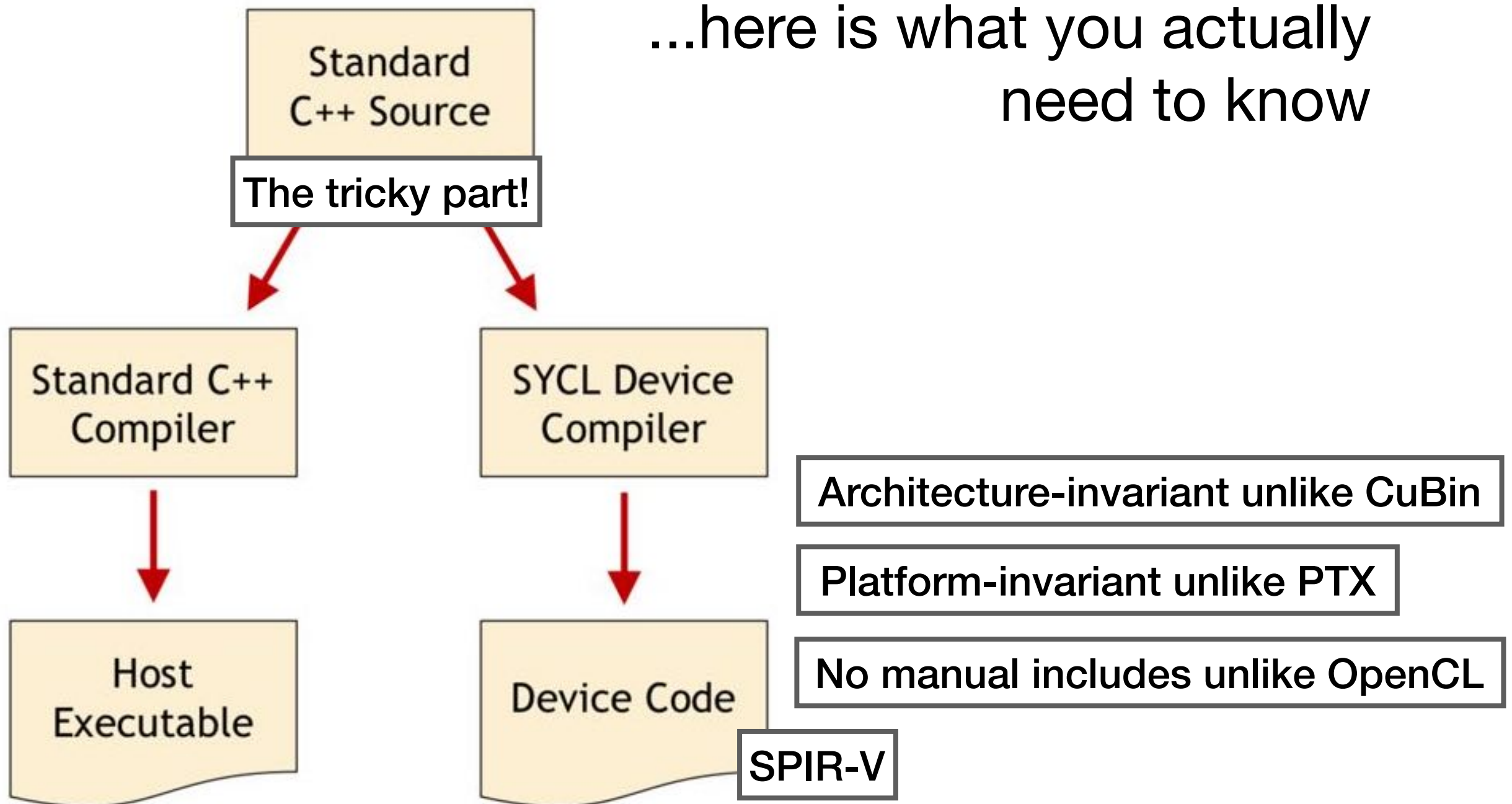


# How SyCL works?



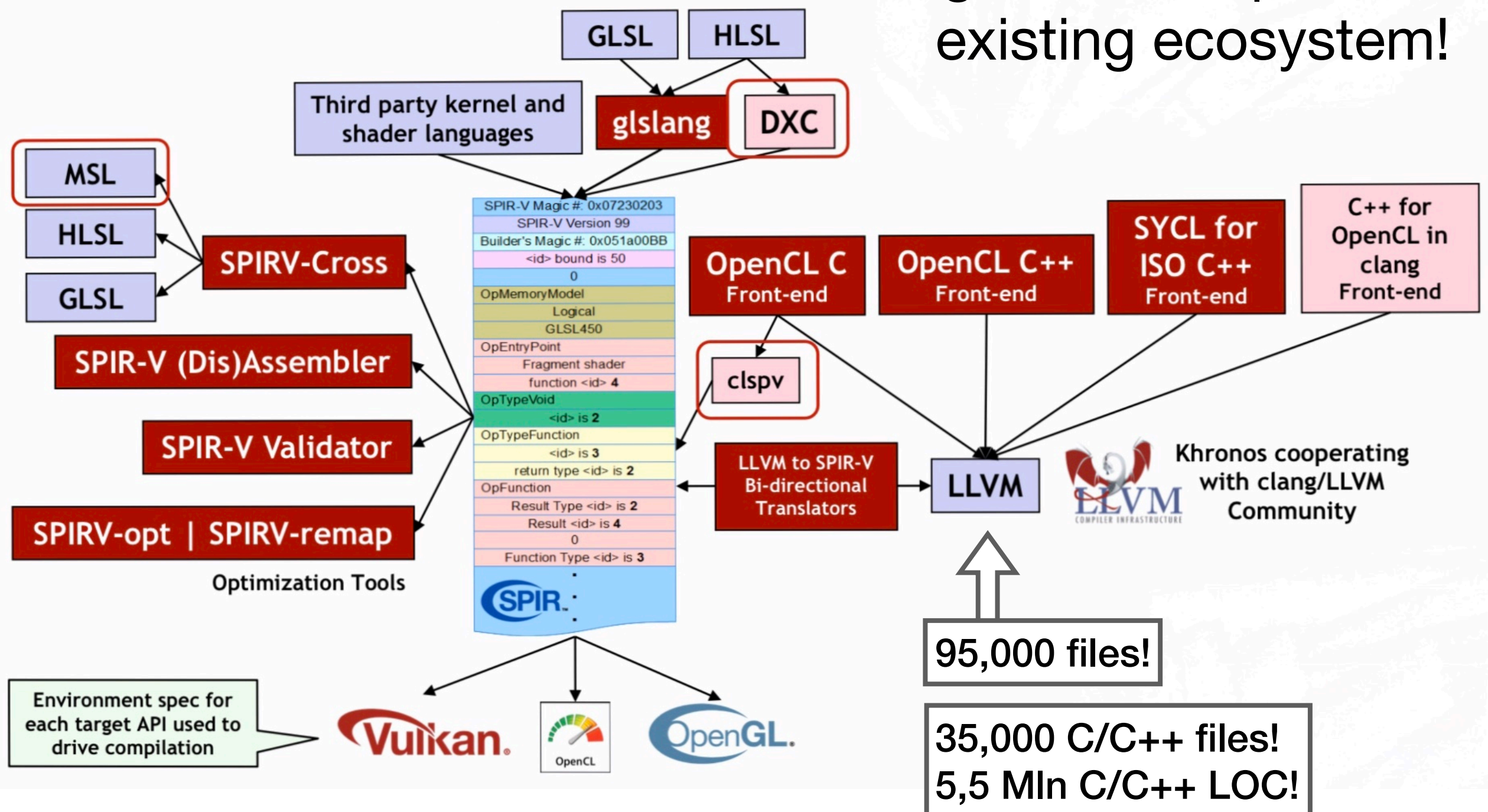
# How SyCL works?

...here is what you actually need to know



# How SyCL works?

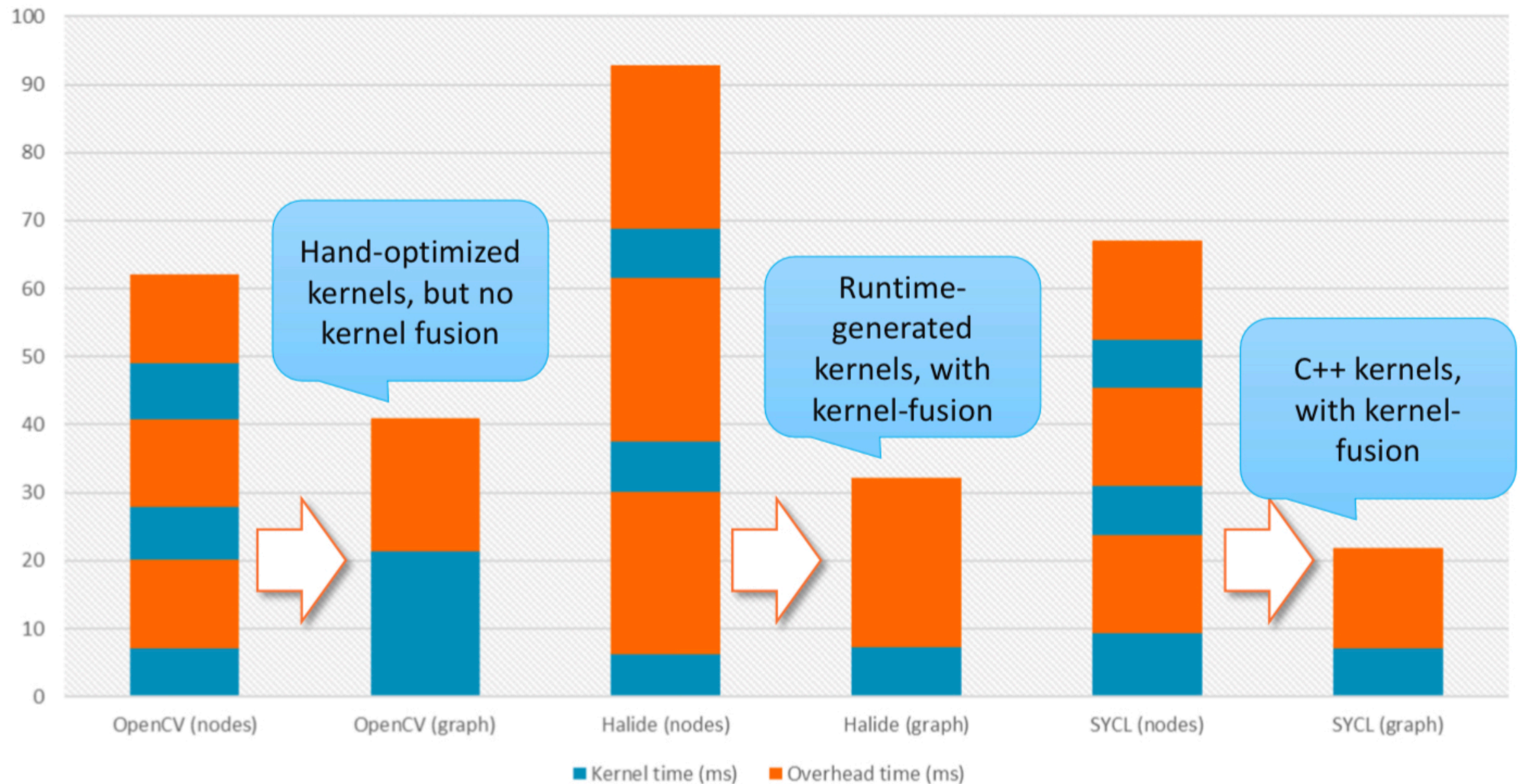
...it grows on top of the existing ecosystem!





# Kernel Fusion

...impact on pipeline performance!

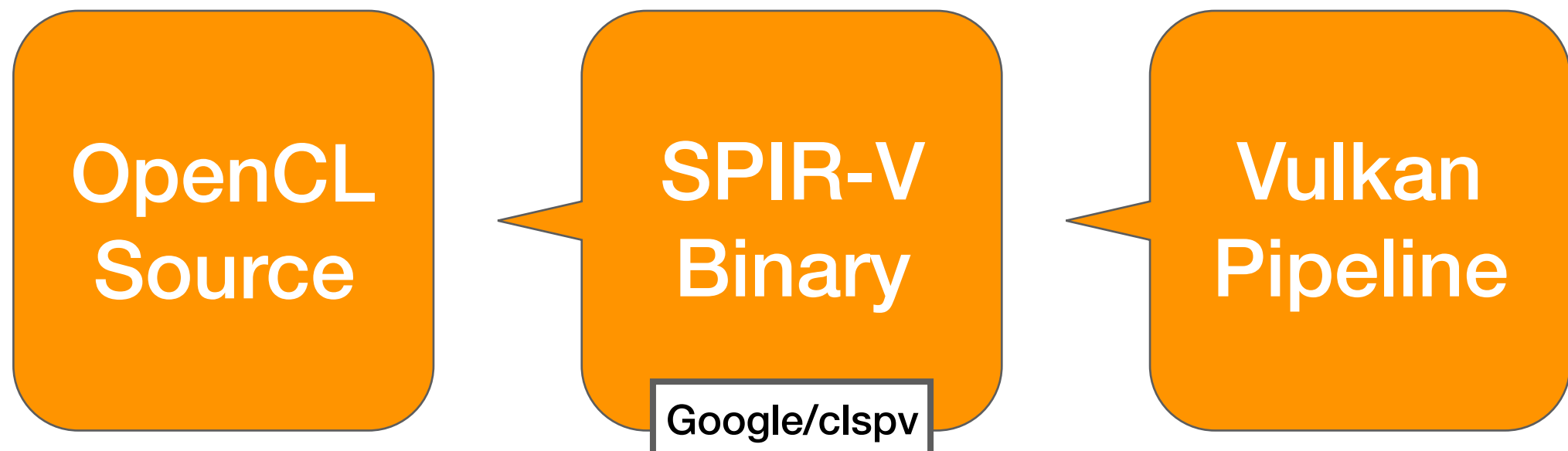


# What to choose?

Compromises

# Unified solution

...if you want cross platform binaries for your custom hand-made kernels!



**OpenCL  
Source**

**SPIR-V  
Binary**

**Vulkan  
Pipeline**

### Pros

**Same binary runs everywhere,  
easy to debug**

**Concurrent queues**

**Logical devices can represent  
SLI groups**

**Same ecosystem for both  
graphics and compute**

### Cons

**Separate CL/C++ codebases**

**Experimental stage compilers**

No support for CUDA features:  
warp shuffles,  
hardware-specific instructions \*,  
L1 cache tuning

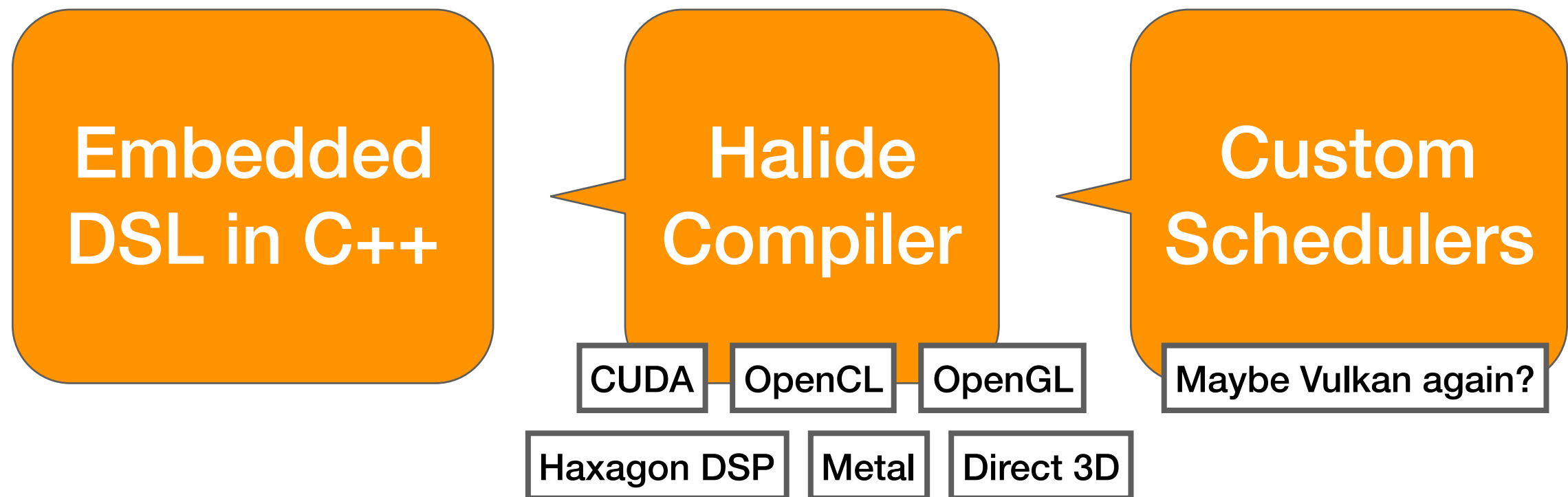
**No work-  
arounds?**

**100 TFlops?!**

No modern C++ support until  
version 2.2

# Flexible solution

...if you want a flexible tool  
here and now!





**Embedded  
DSL in C++**

**Halide  
Compiler**

**Custom  
Scheduler**

### Pros

**Great for prototyping and  
benchmarking**

**Generate binaries for every  
platform**

**Easy to export computational  
graphs into other libs**

Easy to debug algorithms

Built-in tools for image processing

### Cons

**Turing-incomplete**

**Limited number of supported types**

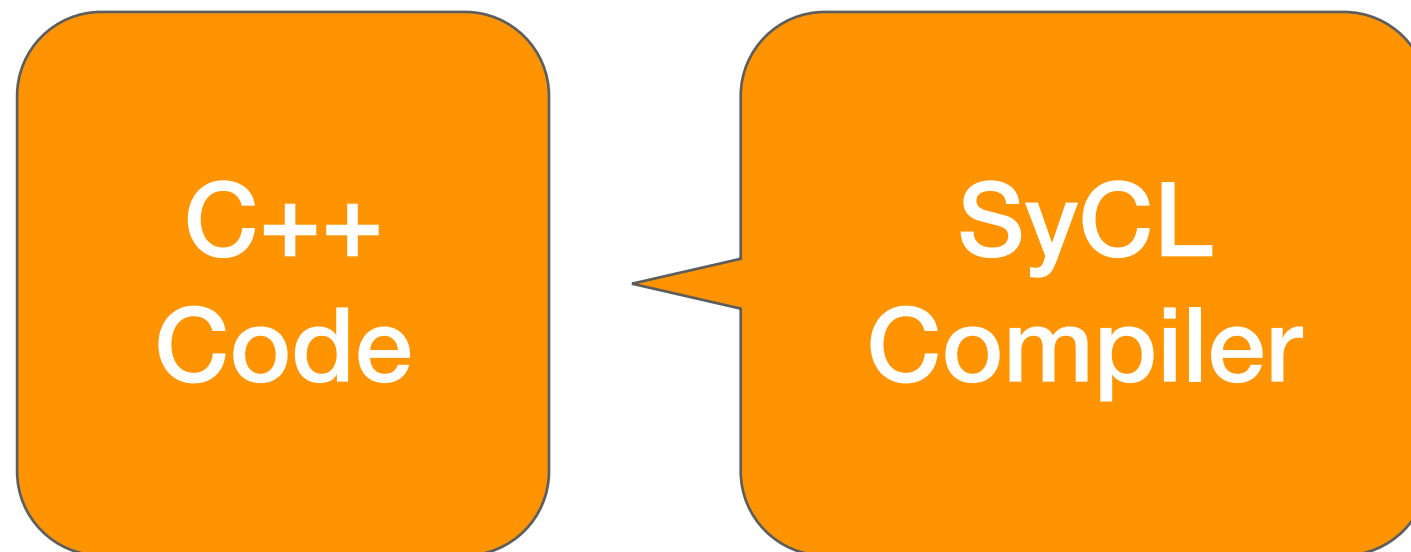
Huge LLVM dependency

Non-standard C++

Bad error messages

# Clean solution

...if you want some classical  
type-safe C++!



**C++  
Code**

**SyCL  
Compiler**

### Pros

**Use C++ templates and lambda functions for host & device code - just pass "sycl" policy**

SYCL will not create C++ language extensions, but instead add features via C++ library

Does kernel fusion

Layered over OpenCL

### Cons

**Very immature, stability and C++17 adoption is expected closer to 2020**

**Underlying implementation requires compiler support**

Kernel fusion may be weak

Boost.Compute dependency

# Simple solution

...if you want a brute-force accelerator for simple data-parallel number-crunching!

High level GPGPU library  
of choice like ArrayFire

# High level GPGPU library of choice like ArrayFire

## Pros

**Already packed with binaries for  
multiple backends**

---

**Minimal coding required**

## Cons

Weak kernel fusion

# Comparison of recipes

...lets summarize our results!

	Simple	Unified	Flexible	Clean
Technology	ArrayFire	CL & SPIR-V	Halide	C++ SyCL
Write code once	Yes	Separate Files	T-Incomplete	Yes
Run everywhere	Almost	Yes	Yes	Eventually
Max performance	Average	High	Highest	High
Minimal code size	Smallest	Mid-Large	Mid-Large	Small

# Comparison of recipes

...lets summarize our results!

	Max Performance Today			
Technology		CL & SPIR-V	Halide	
Write code once		Yes	T-Incomplete	
Run everywhere		Yes	Yes	
Max performance		High	Highest	
Minimal code size		Mid-Large	Mid-Large	

# Comparison of recipes

...lets summarize our results!

	Sometime in the Future			
Technology				C++ SyCL
Write code once				Yes
Run everywhere				Eventually
Max performance				High
Minimal code size				Small





# Bonus: Crazy solution

...if only you are as obsessed with parallel computing as I am!

A New Language

Compiler

Parallelism

Reflections

Simplicity of parsing

For vectorization analysis

Context-free grammars for fast JIT

For serialization and data exchange

# Bonus: Hierarchy

...approximation of high-performance C++  
GPGPU solutions

Symbolic Graph	<b>TF, PyTorch, cuDNN</b> , MKL-DNN
Lazy Evaluation	Eigen, <b>TF</b> , VexCL, ArrayFire
Linear Algebra	Eigen, MKL, VexCL, <b>cuBLAS, ArrayFire</b> , Boost.Compute
Scheduling	Intel TBB, Vulkan, <b>OpenMP</b> , SyCL
Language & Extensions*	<b>CUDA, OpenCL</b> , GLSL, OpenMP, OpenACC
Compilers*	<b>LLVM, TVM</b> , GCC

# Q & A

[github.com/ashvardanian/SandboxGPUs](https://github.com/ashvardanian/SandboxGPUs)

[github.com/ashvardanian](https://github.com/ashvardanian)

[fb.com/ashvardanian](https://fb.com/ashvardanian)

[linkedin.com/in/ashvardanian](https://linkedin.com/in/ashvardanian)

[vk.com/ashvardanian](https://vk.com/ashvardanian)