

Implementing the SYCL for OpenCL Shared Source C++ Programming Model using Clang/LLVM

Gordon Brown
Runtime Engineer, Codeplay

Visit us at
www.codeplay.com

45 York Place
Edinburgh
EH1 3HP
United Kingdom



Agenda

- **Overview of SYCL**
- **SYCL Example: Vector Add**
- **Shared Source Programming Model**
- **Implementing SYCL Using Clang/LLVM**

Overview of SYCL

Motivation

- To make GPGPU simpler and more accessible.
- To create a C++ for OpenCL™ ecosystem.
 - Combine the ease of use and flexibility of C++ and the portability and efficiency of OpenCL.
- To provide a foundation for constructing complex and reusable template algorithms:
 - `parallel_reduce()`, `parallel_map()`, `parallel_sort()`
- To define an open and portable standard.

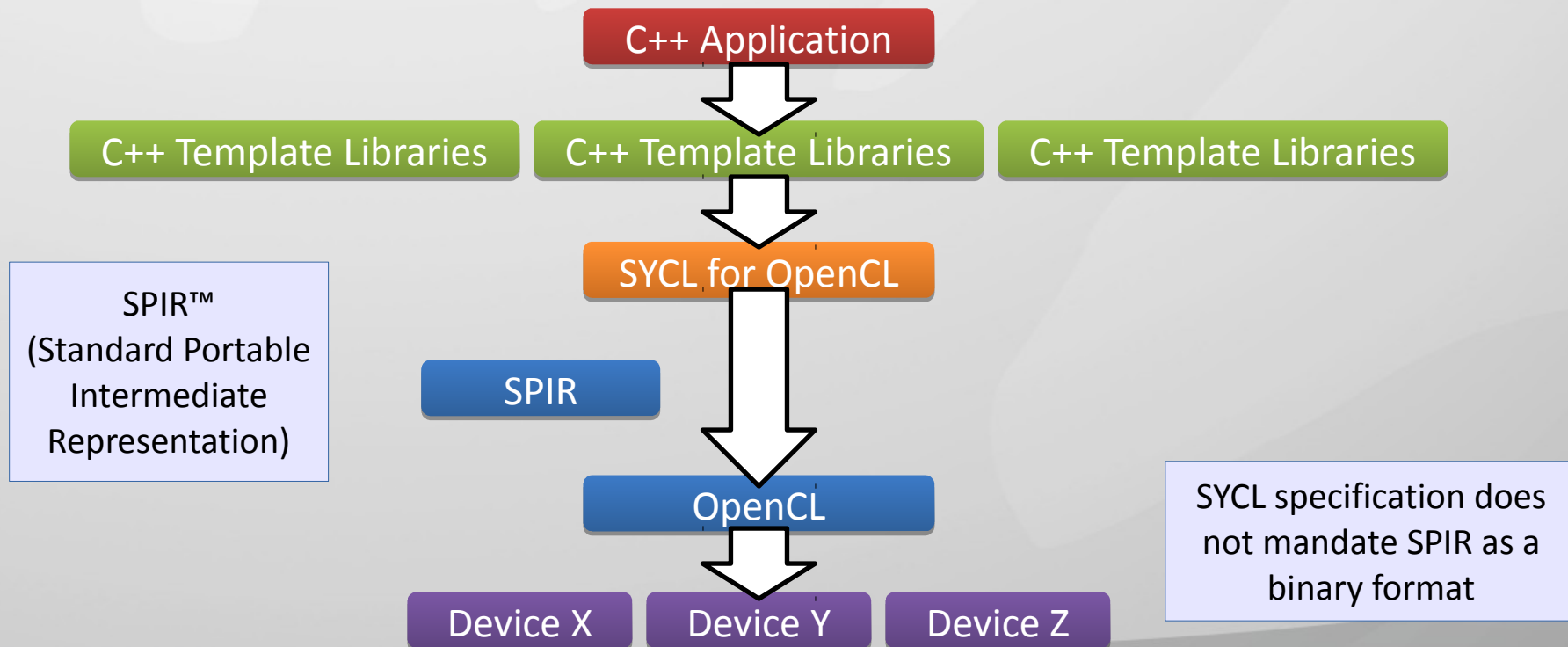
SYCL for OpenCL



Cross platform, single source, C++ programming layer

Built on top of OpenCL and based on standard C++11.

The SYCL Ecosystem



SYCL Standard Roadmap

- Current State:
 - Second provisional specification being announced here at Supercomputing 2014.
- Next Steps:
 - Full specification based on feedback.
 - Conformance test suite to ensure compatibility.
 - Release implementations.

SYCL Example: Vector Add

```
#include <CL/sycl.hpp>
using namespace cl::sycl;
template <typename T>
void parallel_vadd(std::vector<T> inputA, std::vector<T> inputB, std::vector<T> output) {
    buffer<float, 1> inputABuf(inputA, inputA.size());
    buffer<float, 1> inputBBuf(inputB, inputB.size());
    buffer<float, 1> outputBuf(output, output.size());
    queue defaultQueue;
    command_group(defaultQueue, [&] () {
        auto inputAPtr = inputABuf.get_access<access::read>();
        auto inputBPtr = inputBBuf.get_access<access::read>();
        auto outputPtr = outputBuf.get_access<access::write>();
        parallel_for< vadd<T> >(range<1>(output.size()), [=](id<1> idx) {
            ouptPtr[idx] = inputAPtr[idx] + inputBPtr[idx];
        }));
    });
}
```

```
#include <CL/sycl.hpp>  
using namespace cl::sycl;
```

```
int main() {  
    return 0;  
}
```

The SYCL runtime is in sycl.hpp and within the cl::sycl namespace

```
#include <CL/sycl.hpp>
using namespace cl::sycl;

int main() {
    int count = 1024;

    std::vector<float> inputA(count) = { /* input a */ };
    std::vector<float> inputB(count) = { /* input b */ };
    std::vector<float> output(count) = { /* output */ };

    return 0;
}
```

Construct and initialise three `std::vector` objects of 1024 float elements, two inputs and one output.

```
#include <CL/sycl.hpp>
using namespace cl::sycl;
```

```
int main() {
    int count = 1024;
    std::vector<float> inputA(count) = { /* input a */ };
    std::vector<float> inputB(count) = { /* input b */ };
    std::vector<float> output(count) = { /* output */ };
    {
        buffer<float, 1> inputABuf(inputA.data(), inputA.size());
        buffer<float, 1> inputBBuf(inputB.data(), inputB.size());
        buffer<float, 1> outputBuf(output.data(), output.size());
    }
    return 0;
}
```

Construct three SYCL buffers and initialise them with the data from the std::vectors.

Data is synchronised by RAI

```
.#include <CL/sycl.hpp>
using namespace cl::sycl;

int main() {
    int count = 1024;
    std::vector<float> inputA(count) = { /* input a */ };
    std::vector<float> inputB(count) = { /* input b */ };
    std::vector<float> output(count) = { /* output */ };
    {
        buffer<float, 1> inputABuf(inputA.data(), inputA.size());
        buffer<float, 1> inputBBuf(inputB.data(), inputB.size());
        buffer<float, 1> outputBuf(output.data(), output.size());
        queue defaultQueue;
    }
    return 0;
}
```

There are many other options for device discovery and configuration.

Construct a SYCL queue to execute work on a device.

...

```
int main() {  
    int count = 1024;  
    std::vector<float> inputA(count) = { /* input a */ };  
    std::vector<float> inputB(count) = { /* input b*/ };  
    std::vector<float> output(count) = { /* output */ };  
  
    {  
        buffer<float, 1> inputABuf(inputA.data(), inputA.size());  
        buffer<float, 1> inputBBuf(inputB.data(), inputB.size());  
        buffer<float, 1> outputBuf(output.data(), output.size());  
        queue defaultQueue;  
        command_group(defaultQueue, [&] () {  
        });  
    }  
    return 0;  
}
```

The `command_group` is en-queued asynchronously and is thread safe.

Construct a SYCL `command_group` to define the work to be en-queued on a device.

```
...  
std::vector<float> inputA(count) = { /* input a */ };  
std::vector<float> inputB(count) = { /* input b*/ };  
std::vector<float> output(count) = { /* output */ };  
{  
    buffer<float, 1> inputABuf(inputA.data(), inputA.size());  
    buffer<float, 1> inputBBuf(inputB.data(), inputB.size());  
    buffer<float, 1> outputBuf(output.data(), output.size());  
    queue defaultQueue;  
    command_group(defaultQueue, [&] () {  
        auto inputAPtr = inputABuf.get_access<access::read>();  
        auto inputBPtr = inputBBuf.get_access<access::read>();  
        auto outputPtr = outputBuf.get_access<access::write>();  
    });  
}  
...
```

The SYCL runtime used accessors to track dependencies across command_groups.

Construct three SYCL accessors with the appropriate access modes, to give the device access to the data.


```
...  
{  
    buffer<float, 1> inputABuf(inputA.data(), inputA.size());  
    buffer<float, 1> inputBBuf(inputB.data(), inputB.size());  
    buffer<float, 1> outputBuf(output.data(), output.size());  
    queue defaultQueue;  
    command_group(defaultQueue, [&] () {  
        auto inputAPtr = inputABuf.get_access<access::read>();  
        auto inputBPtr = inputBBuf.get_access<access::read>();  
        auto outputPtr = outputBuf.get_access<access::write>();  
        parallel_for<class vadd>(range<1>(count), ([=](id<1> idx) {  
            }));  
    });  
}  
...
```

There are additional more complex APIs.

Call `parallel_for()` to execute a kernel function.

The typename 'vadd' is used to name the lambda.

The range provided to the `parallel_for()` should match the size of the data buffers.

```
...  
{  
    buffer<float, 1> inputABuf(inputA.data(), inputA.size());  
    buffer<float, 1> inputBBuf(inputB.data(), inputB.size());  
    buffer<float, 1> outputBuf(output.data(), output.size());  
    queue defaultQueue;  
    command_group(defaultQueue, [&] () {  
        auto inputAPtr = inputABuf.get_access<access::read>();  
        auto inputBPtr = inputBBuf.get_access<access::read>();  
        auto outputPtr = outputBuf.get_access<access::write>();  
        parallel_for<class vadd>(range<1>(count), ([=](id<1> idx) {  
            outputPtr[idx] = inputAPtr[idx] + inputBPtr[idx];  
        }));  
    });  
}  
...
```

The body of the lambda expression is what is compiled into an OpenCL kernel by the SYCL device compiler.

Use the subscript operator on the accessors to read and write the data.

```
#include <CL/sycl.hpp>
```

```
using namespace cl::sycl;
```

```
void parallel_vadd(std::vector<float> &inputA, std::vector<float> &inputB, std::vector<float> &output) {  
    buffer<float, 1> inputABuf(inputA.data(), inputA.size());  
    buffer<float, 1> inputBBuf(inputB.data(), inputB.size());  
    buffer<float, 1> outputBuf(output.data(), output.size());  
    queue defaultQueue;  
    command_group(defaultQueue, [&] () {  
        auto inputAPtr = inputABuf.get_access<access::read>();  
        auto inputBPtr = inputBBuf.get_access<access::read>();  
        auto outputPtr = outputBuf.get_access<access::write>();  
        parallel_for<class vadd>(range<1>(count), ([=](id<1> idx) {  
            outputPtr[idx] = inputAPtr[idx] + inputBPtr[idx];  
        }));  
    });  
}
```

Create a function that takes the input and output vectors

```
#include <CL/sycl.hpp>
using namespace cl::sycl;
template <typename T>
void parallel_vadd(std::vector<T> &inputA, std::vector<T> &inputB, std::vector<T> &output) {
    buffer<T, 1> inputABuf(inputA.data(), inputA.size());
    buffer<T, 1> inputBBuf(inputB.data(), inputB.size());
    buffer<T, 1> outputBuf(output.data(), output.size());
    queue defaultQueue;
    command_group(defaultQueue, [&] () {
        auto inputAPtr = inputABuf.get_access<access::read>();
        auto inputBPtr = inputBBuf.get_access<access::read>();
        auto outputPtr = outputBuf.get_access<access::write>();
        parallel_for<vadd<T>>(range<1>(count), ([=](id<1> idx) {
            outputPtr[idx] = inputAPtr[idx] + inputBPtr[idx];
        }));
    });
}
```

Template the function
By the data type

The typename 'vadd' must
also be templated as the lambda
expression is template dependant.

Comparison with OpenCL

```

#include <stdio.h>
#include <stdlib.h>
#include <sys/types.h>
#include <string.h>
#include <OpenCL/opencl.h>
#include <vector>
#include <string>
#include <CL/sycl>
using namespace cl;

//[3] Set up device type from compiler command line or from
//the default type
sycl::device DEVICE;
#define DEVICE_TYPE_DEVICE sycl::device
#define DEVICE_TYPE_HOST sycl::host

//[4] Set up device type from compiler command line or from
//the default type
sycl::device DEVICE;
#define DEVICE_TYPE_DEVICE sycl::device
#define DEVICE_TYPE_HOST sycl::host

int main(int argc, char** argv)
{
    // Get the device type from the compiler command line or from the default type
    sycl::device dev;
    if (argc > 1) dev = sycl::device(argv[1]);
    else dev = sycl::device();

    // Create a compute queue
    sycl::queue q(dev);

    // Create a command queue
    sycl::command_queue cq(dev, q);

    // Create the program from the source buffer
    sycl::program prog(argv[2], argv[3], argv[4], argv[5], argv[6], argv[7], argv[8], argv[9], argv[10], argv[11], argv[12], argv[13], argv[14], argv[15], argv[16], argv[17], argv[18], argv[19], argv[20], argv[21], argv[22], argv[23], argv[24], argv[25], argv[26], argv[27], argv[28], argv[29], argv[30], argv[31], argv[32], argv[33], argv[34], argv[35], argv[36], argv[37], argv[38], argv[39], argv[40], argv[41], argv[42], argv[43], argv[44], argv[45], argv[46], argv[47], argv[48], argv[49], argv[50], argv[51], argv[52], argv[53], argv[54], argv[55], argv[56], argv[57], argv[58], argv[59], argv[60], argv[61], argv[62], argv[63], argv[64], argv[65], argv[66], argv[67], argv[68], argv[69], argv[70], argv[71], argv[72], argv[73], argv[74], argv[75], argv[76], argv[77], argv[78], argv[79], argv[80], argv[81], argv[82], argv[83], argv[84], argv[85], argv[86], argv[87], argv[88], argv[89], argv[90], argv[91], argv[92], argv[93], argv[94], argv[95], argv[96], argv[97], argv[98], argv[99]);

    // Build the program
    prog.build(q);

    // Create the kernel
    sycl::kernel k(prog, "kernel");

    // Create the input arrays in device memory
    sycl::buffer<int> a(q, {0, 0, 0});
    sycl::buffer<int> b(q, {0, 0, 0});
    sycl::buffer<int> c(q, {0, 0, 0});

    // Set the arguments to our compute kernel
    auto args = {a, b, c};

    // Set the arguments to our compute kernel
    auto args = {a, b, c};

    // Execute the kernel
    k.execute(args);

    // Get the output arrays
    a.get_host(&out_a);
    b.get_host(&out_b);
    c.get_host(&out_c);

    // Print the results
    for (int i = 0; i < out_a.size(); i++)
        printf("%d\t", out_a[i]);
    for (int i = 0; i < out_b.size(); i++)
        printf("%d\t", out_b[i]);
    for (int i = 0; i < out_c.size(); i++)
        printf("%d\t", out_c[i]);

    return 0;
}
    
```

```

#include <CL/sycl.hpp>
using namespace cl::sycl;

int main()
{
    int count = 1024;

    sycl::queue q(sycl::device());

    sycl::buffer<int> a(sycl::range<1>(count));
    sycl::buffer<int> b(sycl::range<1>(count));
    sycl::buffer<int> c(sycl::range<1>(count));

    auto p0 = a.get_access(sycl::read);
    auto p1 = b.get_access(sycl::read);
    auto p2 = c.get_access(sycl::write);

    auto kernel = sycl::kernel{q, "kernel", {0, 0, 0}};

    kernel.execute(a, b, c);

    return 0;
}
    
```

SYCL for OpenCL

```

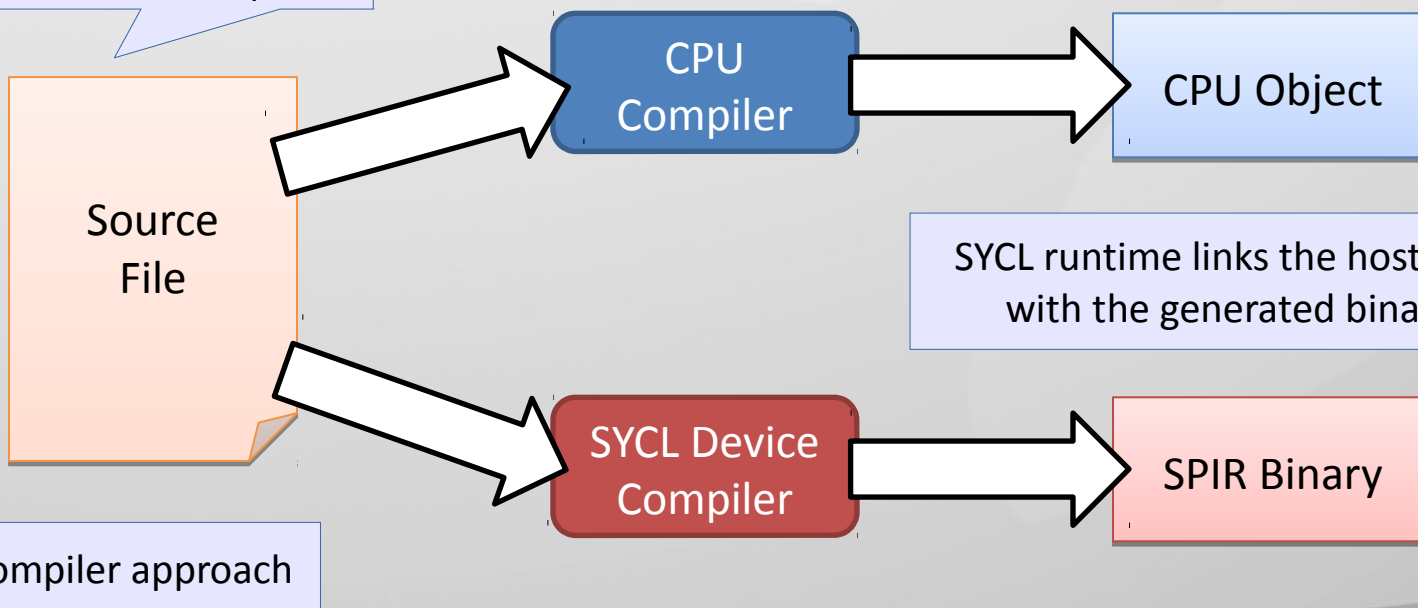
...kernel void void
..._global float* a
..._global float* b
..._global float* c
const unsigned int count;
int i = 0;
for (i = 0; i < count; i++)
    c[i] = a[i] + b[i];
}
    
```

Traditional OpenCL

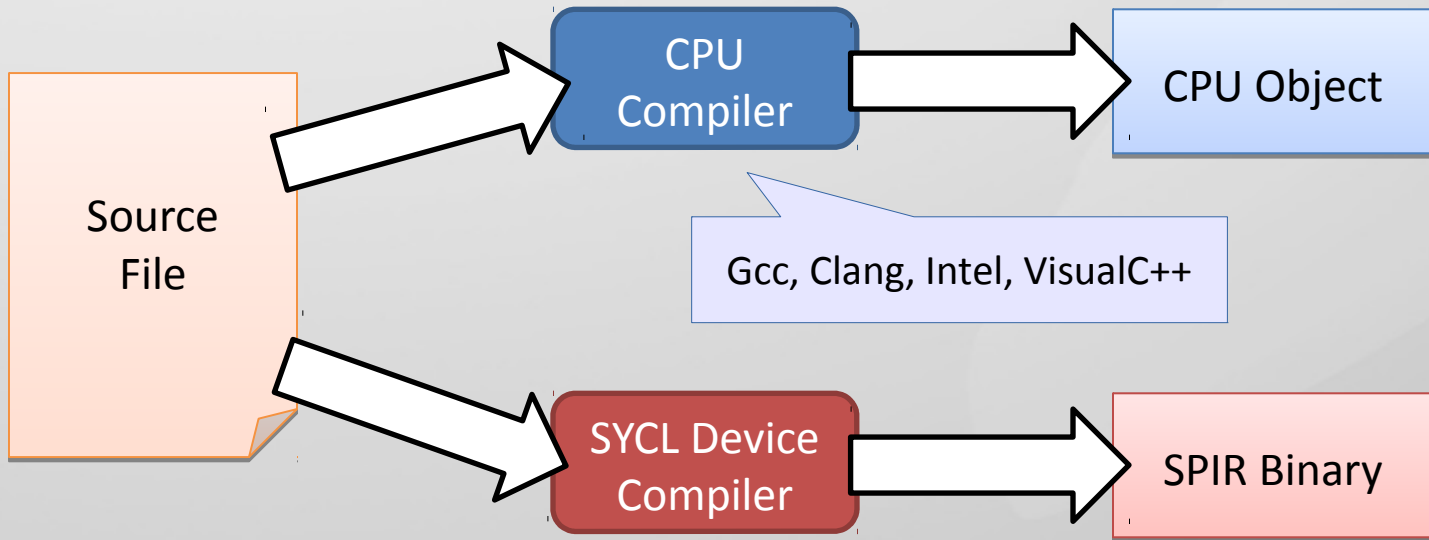
Shared Source Programming Model

Shared Source

Single source file compiled by host and device compilers

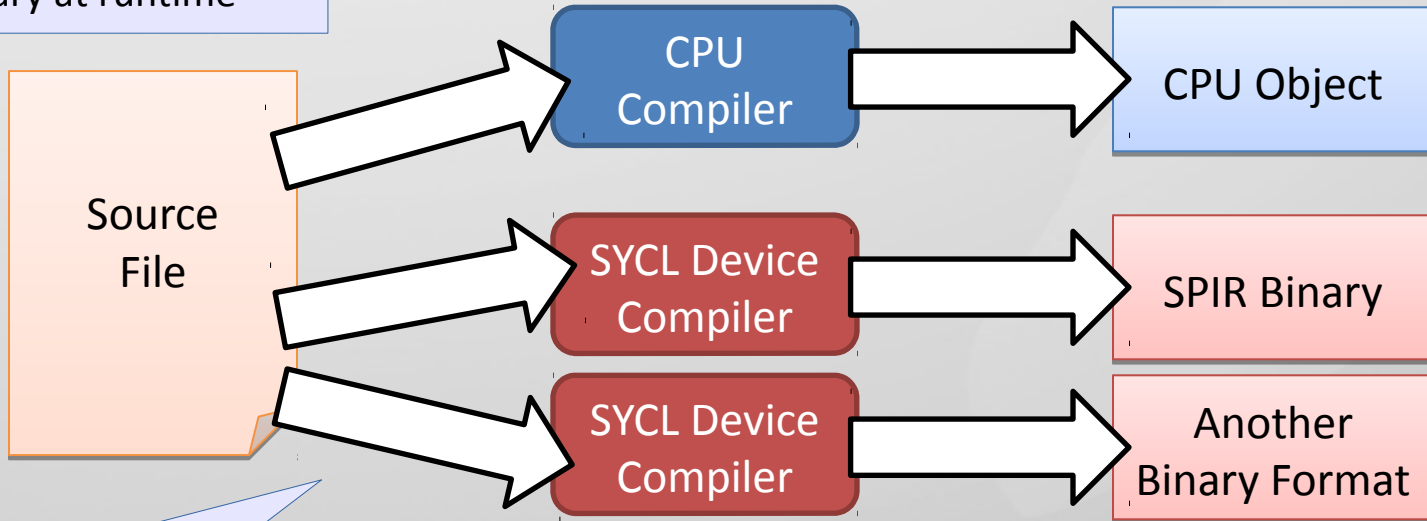


Choose Your Own Host Compiler



Target Multiple Device Compilers

SYCL runtime will choose
a binary at runtime



Multiple device compilers

Host Device

Can be used as a target for debugging

SYCL Runtime

Can be used as an optional fall back if a Kernel execution fails

OpenCL

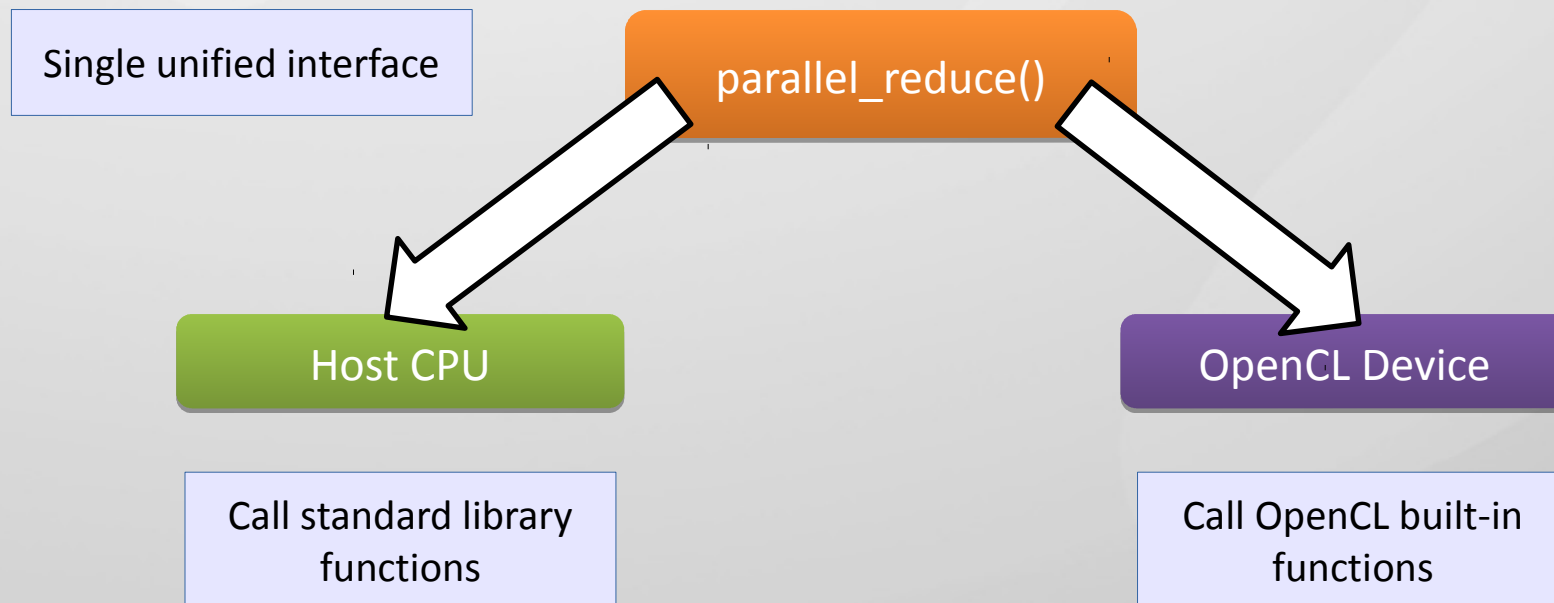
Host Fallback

Device X

Device Y

Device Z

Use Common Libraries



Implementing SYCL Using Clang/LLVM

Why Clang/LLVM?

- Clang & LLVM is perfect for implementing innovative compiler technologies:
 - Large amount of contributors and a great developer community.
 - Very feature rich (standard and non-standard).
 - Full of re-usable modules and components.

Topics

- Separating Host & Device Code
- Constructing an OpenCL Kernel Function
- Duplicating Device Functions
- Diagnosing Invalid Device Code
- Supporting OpenCL Types
- Generating a SPIR Binary
- Integrating with SYCL Runtime

Separating Host & Device Code

SYCL Requirement

- SYCL has a unified interface that users develop in.
- Programming model needs to allow the host compiler and SYCL device compiler to extract different code:
 - Types that can be passed between host and device.
 - APIs that have different behavior between host and device.

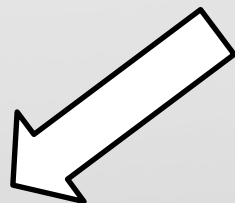
Implicit Pre-process Macro

- Our SYCL device compiler adds the pre-process macro.
- This is used to separate parts of the SYCL runtime into host and device code.

```
__SYCL_DEVICE_ONLY__
```

Example: accessor

```
accessor<float, 1, access::read> inputPtrA(inputBufferA);
```



Host Compiler

Construct OpenCL buffer
and enqueue to device



Device Compiler

Represent OpenCL kernel
argument on device

Example: parallel_for()

```
parallel_for<class vadd>(range<1>(count), ([=](id<1> itemID) {  
    outputPtr[itemID] = inputPtrA[itemID] + inputPtrB[itemID];  
}));
```



Host Compiler

Set arguments and
enqueue kernel function

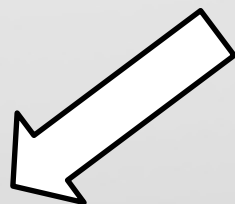


Device Compiler

Generate the
kernel function

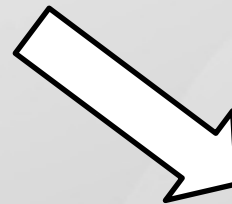
Example: fmin()

```
float4 result = fmin(valA, valB);
```



Host Compiler

Host-side fmin()
function



Device Compiler

OpenCL built-in fmin()
function

Constructing an OpenCL Kernel Function

SYCL Requirement

- SYCL represents kernel functions using lambda expressions.
- Programming model needs to identify and construct an OpenCL kernel:
 - Name the lambda expression.
 - Convert the body to an OpenCL kernel function.
 - Convert captured variables to OpenCL kernel arguments.

Identify the SYCL Kernel Function

- Our SYCL device compiler uses a Clang attribute to identify the SYCL kernel function.
- This is also used to assign a name to the lambda expression.

```
__attribute__((sycl_kernel(name)))
```

Convert the Lambda Expression

- Our SYCL device compiler manipulates the Clang AST to create the OpenCL kernel function.
 - The lambda expression is converted to an OpenCL compatible function.
 - The captured variables are converted to OpenCL kernel arguments.

Example: parallel_for()

```
parallel_for<class vadd>(range<1>(count), ([=](id<1> itemID) {  
    outputPtr[itemID] = inputPtrA[itemID] + inputPtrB[itemID];  
}));
```



```
__kernel void sycl_vadd (__global float *inputPtr A,  
                        __global float *inputPtrB,  
                        __global float *outputPtr,  
                        id itemID) {  
    outputPtr[itemID] = inputPtrA[itemID] + inputPtrB[itemID];  
}
```

Duplicating Device Functions

SYCL Requirement

- SYCL requires all functions that are called by the SYCL kernel function to also be available.
- Programming model needs to transform all called functions into device compatible functions:
 - Traverse call graph.
 - Clone function definitions.
 - Deduce OpenCL address spaces for parameters.

Transforming Device Functions

- Out SYCL device compiler uses the Clang TreeTransform framework:
 - Traverses the Clang AST.
 - Identifies the OpenCL kernel function's call graph.
 - Transforms called functions into device functions.
 - Deduces address spaces for device function based on parameters.

Example: reduction

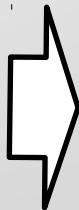
```
template <typename N, typename L>  
void parallel_for(range<1>, L);
```



```
template <typename T>  
T reduce_array(T *data);
```



```
template <typename T>  
T reduce(T a, T b);
```



```
__kernel void sycl_kernel0(__global float *ptr);
```



```
float reduce_array(__global float *data);
```



```
float reduce(float a, float b);
```

Diagnosing Invalid Device Code

SYCL Requirement

- There are restrictions on the C++ features that SYCL can support based on OpenCL 1.x hardware support.
- Programming model needs to diagnose device code when invalid C++ features are being used:
 - Catch device code that contains invalid features.
 - Generate useful error messages.

Unsupported C++ Features

- Recursion.
- Exception Handling.
- RTTI.
- Dynamic Allocation.
- Dynamic Polymorphism.
- Static Variables.
- Function Pointers.
- Virtual Functions.

These apply only to SYCL kernel functions

AST Checkers

- Our SYCL device compiler performs device code diagnostics using AST checkers.
- These are constructed from Clang's AST matcher framework.

Example: exception handling checker

```
struct ExceptionChecker : internal::SYCLFunctionBodyChecker<Stmt> {  
    MatcherType<Stmt> getMatcher() const override {  
        return stmt(anyOf(tryStmt(), throwExpr()));  
    }  
    void checkMatch(const Stmt *S, const BoundNodes &) override {  
        diagInvalidConstruct(S, "Using exception handling");  
    }  
};
```

Supporting OpenCL Types

SYCL Requirement

- SYCL requires its types to map to underlying OpenCL types transparently and with little or no overhead.
- Programming model needs support all OpenCL types underneath SYCL types:
 - Address spaces.
 - Vectors & swizzles.
 - OpenCL built-in types and qualifiers.

OpenCL Support

- Our SYCL device compiler supports OpenCL types by enabling Clangs built-in support:
 - `__attribute__((address_space(n)))`
 - `__attribute__((ext_vector_type(n)))`
 - `OCLImage2d`, `OCLSampler`, ...
 - `OpenCLImageAccess`

Example: vectors

```
float4 ( a.k.a. vec<float, 4> )
```



```
float __attribute__((ext_vector_type(4)))
```

Example: swizzles

```
float4 lhs, rhs;  
lhs.wzyx() = rhs.xzyw();
```



```
float __attribute__((ext_vector_type(4))) lhs, rhs;  
lhs.wzyx = rhs.xywx;
```

Example: images

```
accessor<float4, 2, access::read, access::image> ptr(img);
```



```
__read_only image2d_t ptr;
```


Example: address spaces

```
accessor<float, 1, access::read, access::global_buffer> ptr(buf);
```



```
__global float *ptr;
```

Generating a SPIR Binary

SYCL Requirement

- SYCL requires the OpenCL kernel function and other functions to be compiled into a binary that can be built from the OpenCL runtime.
- Programming model needs to generate a SPIR output:
 - SPIR binary.
 - SPIR assembly

Output SPIR Bitcode

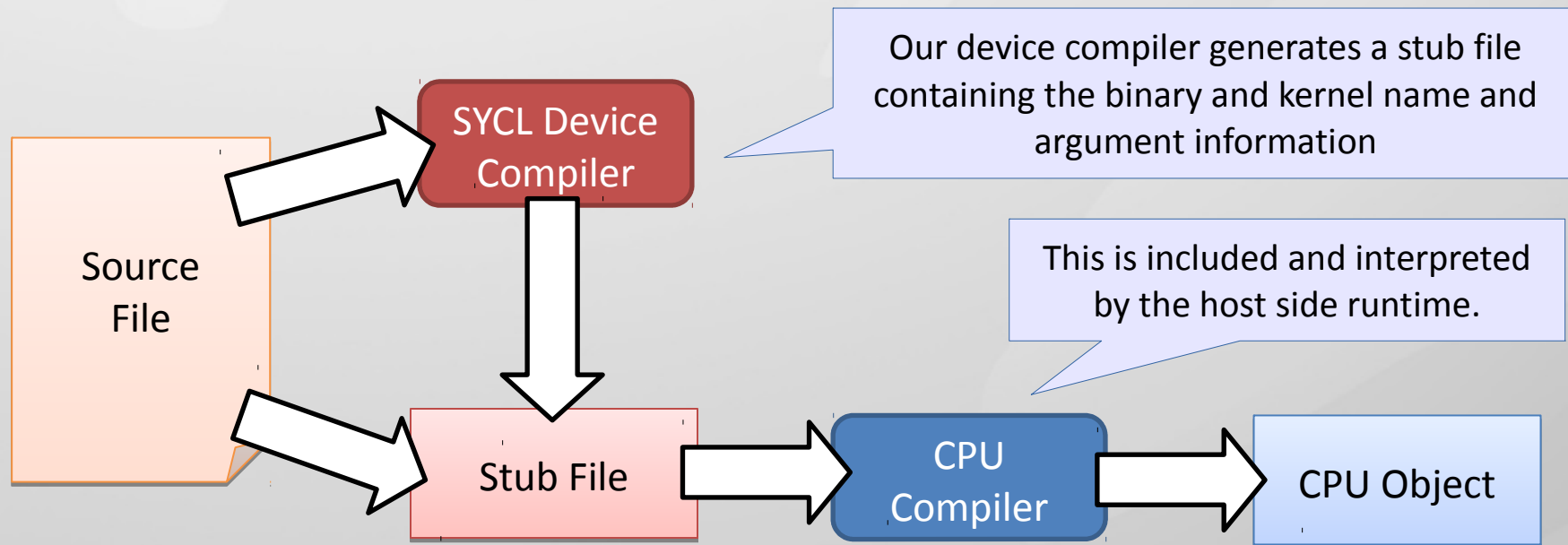
- Our SYCL device compiler outputs a SPIR binary using Clang's support for dumping out SPIR IR.
- It also makes use of Clang's support for generating SPIR meta data.

Integrating with SYCL Runtime

SYCL Requirement

- SYCL requires the runtime to build and execute the SPIR binary.
- Programming model needs to incorporate a way for linking device code into the runtime:
 - OpenCL kernel function binary.
 - OpenCL kernel function name.
 - OpenCL kernel argument information.

Shared Source Integration



How to Get Involved

- OpenCL: Version 2.0 and Beyond
 - Tuesday 5:30PM – 7:00PM
(Rooms 275 – 277)
- SYCL demos:
 - AMD booth in exhibition
- Come and speak to us here:
 - Myself and Codeplay CEO
Andrew Richards
- Read the specification:
 - <https://www.khronos.org/opencl/sycl>
- Leave feedback on the Khronos forum:
 - <https://www.khronos.org/opencl/sycl>
- Try out the open source implementation:
 - <https://www.github.com/amd/triSYCL>