# Extending DPC++ with Support for Huawei Ascend AI Chipset

**Wilson Feng**

**Rasool Maghareh**

**Amy Wang**

Huawei Heterogeneous Compiler Lab, Canada

IWOCL & SYCLcon 2021

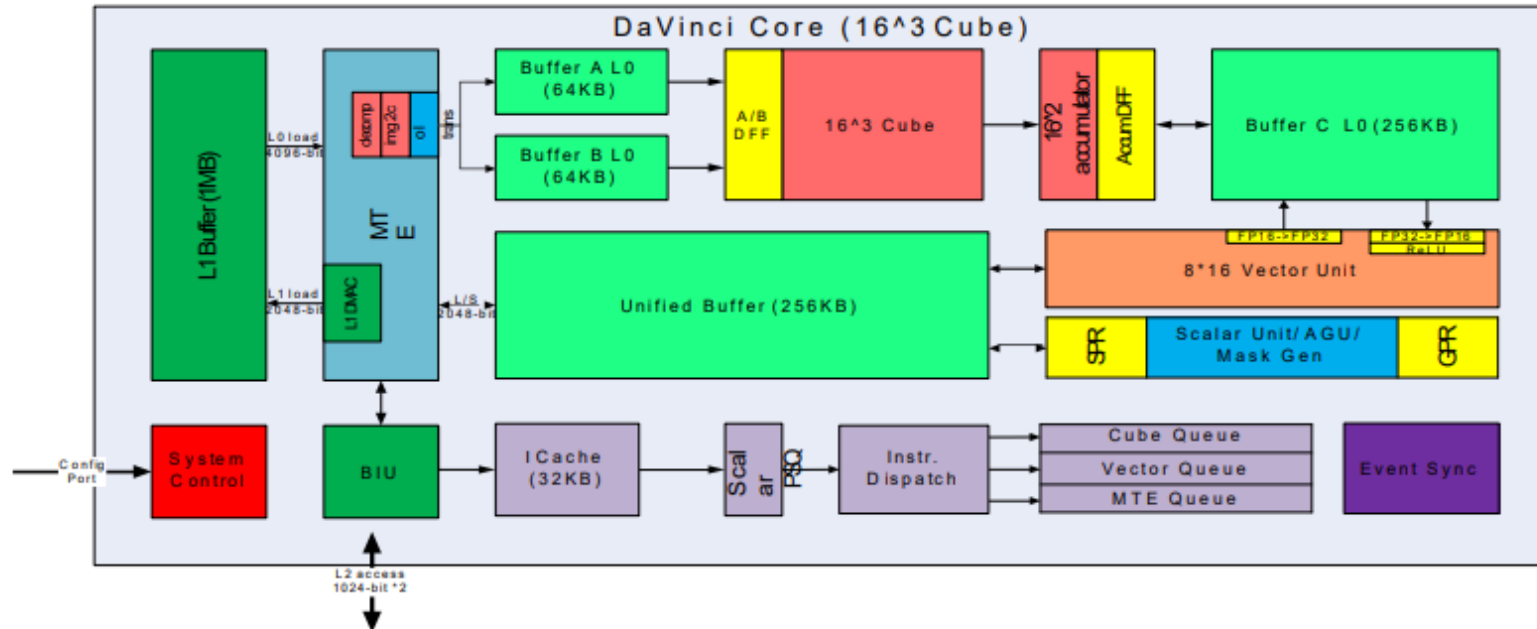April 2021

# Agenda

- Huawei Ascend AI Background

- Our Contribution to DPC++
  - CCE SYCL Backend
  - CCE SYCL Plugin
  - Compilation Toolchain
  - Extension to USM
  - Support for Parallel_for

- Supported Examples

- Future Work

# Huawei Ascend AI Background

- Huawei's custom SoC ASIC for AI workloads

- Host-Device programming model
  - Generic C++ Host Code
  - CCE Device Code – C/C++ based programming language for Ascend AI devices
    - Some C++ features are disabled
    - Explicit software management on different hardware pipelines (DMA transfers and synchronization)
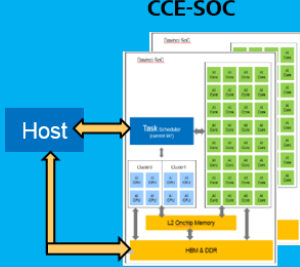
# Huawei Ascend AI Background

**DaVinci Core**



- **Cube**: 4096($16^3$) FP16 MACs + 8192 INT8 MACs
- **Vector**: 2048bit INT8/FP16/FP32 vector with special functions
  (activation functions, NMS- Non Minimum Suppression, ROI, SORT)
- Explicit memory hierarchy design, managed by MTE

# CCEC Compilation



- The source codes are firstly divided into three parts according to the function attribute.
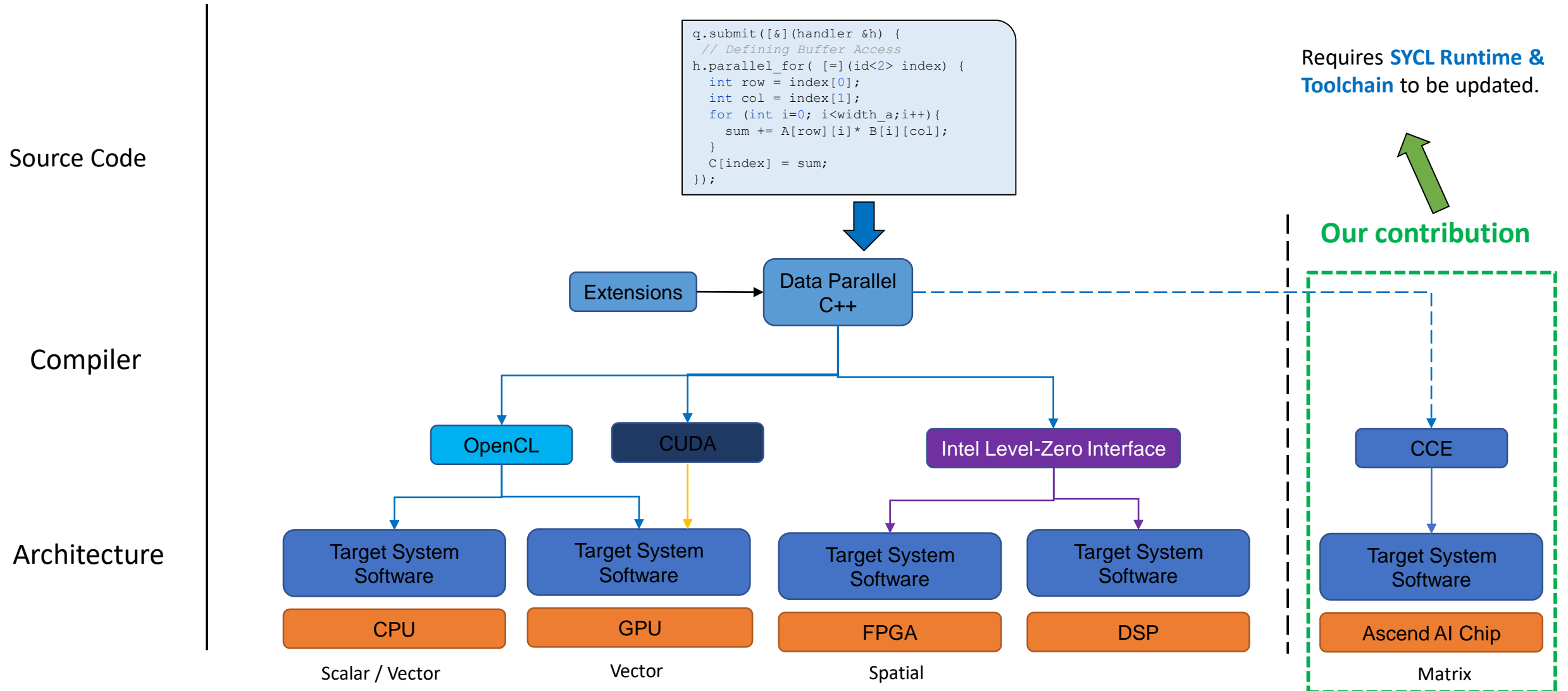- We use different compilers to compile different types of source codes, then we get the **linked** aicore and aicpu binary and relocatable host objects.
- The key step is to link device binaries and host objects together.

# Our Contribution: CCE SYCL Plugin

```
q.submit([&](handler &h) {
  // Defining Buffer Access
  h.parallel_for( [=](id<2> index) {
    int row = index[0];
    int col = index[1];
    for (int i=0; i<width_a;i++){
      sum += A[row][i]* B[i][col];
    }
    C[index] = sum;
  });
```

Requires **SYCL Runtime & Toolchain** to be updated.

**Our contribution**

**Source Code**

**Compiler**

Extensions → Data Parallel C++

OpenCL    CUDA    Intel Level-Zero Interface    CCE

**Architecture**

| Target System Software | Target System Software | Target System Software | Target System Software | Target System Software |
|---|---|---|---|---|
| CPU | GPU | FPGA | DSP | Ascend AI Chip |
| Scalar / Vector | Vector | Spatial | | Matrix |

# Our Contribution: CCE SYCL Backend

**CCE Backend added to SYCL Runtime:**

1. **Device** discovery and selection

2. **Platform** interface

3. **Context** interface

4. **Queue & Event** interfaces

**Selector classes in DPC++:**

- host_selector, cpu_selector, gpu_selector,…

- Added **hiipu_selector()** which selects Huawei Ascend AI Chipset

**cl**::**sycl**::queue Queue(**cl**::**sycl**::**hiipu_selector**{});

sycl-ls --verbose
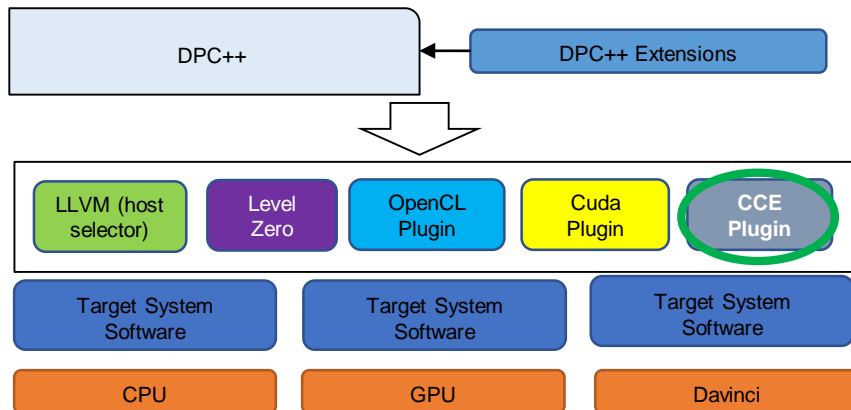


```
Platforms: 0x3
Platform [#0x1]:
    Version  : CUDA 11.1
    Name     : NVIDIA CUDA BACKEND
    Vendor   : NVIDIA Corporation
    Devices  : 0x1
        Device [#0x1]:
            Type      : GPU
            Version   : PI 0.0
            Name      : Tesla P100-PCIE-12GB
            Vendor    : NVIDIA Corporation
            Driver    : CUDA 11.1
Platform [#0x2]:
    Version  : CCE Model_V200
    Name     : Huawei CCE BACKEND
    Vendor   : Huawei Corporation
    Devices  : 0x1
        Device [#0x1]:
            Type      : HiIPU
            Version   : PI 0.0
            Name      : CCE Model_V200
            Vendor    : Huawei Corporation
            Driver    : Model_V200
Platform [#0x3]:
    Version  : 1.2
    Name     : SYCL host platform
    Vendor   :
    Devices  : 0x1
        Device [#0x1]:
            Type      : HOST
            Version   : 1.2
            Name      : SYCL host device
            Vendor    :
            Driver    : 1.2
default_selector()     : GPU : PI 0.0[ CUDA 11.1 ]
host_selector()        : HOST: 1.2[ 1.2 ]
accelerator_selector() : No device of requested type available. -1 (CL_DEVI...
cpu_selector()         : No device of requested type available. -1 (CL_DEVI...
gpu_selector()         : GPU : PI 0.0[ CUDA 11.1 ]
hiipu_selector()       : HiIPU : PI 0.0[ Model_V200 ]
custom_selector(gpu)   : GPU : PI 0.0[ CUDA 11.1 ]
custom_selector(hiipu) : HiIPU : PI 0.0[ Model_V200 ]
custom_selector(cpu)   : No device of requested type available. -1 (CL_DEVI...
custom_selector(acc)   : No device of requested type available. -1 (CL_DEVI...
```

# Our Contribution: CCE SYCL Plugin

**CCE Runtime Plugin:** The runtime plugin performs **command group scope** instructions which act as an interface between the host and device.
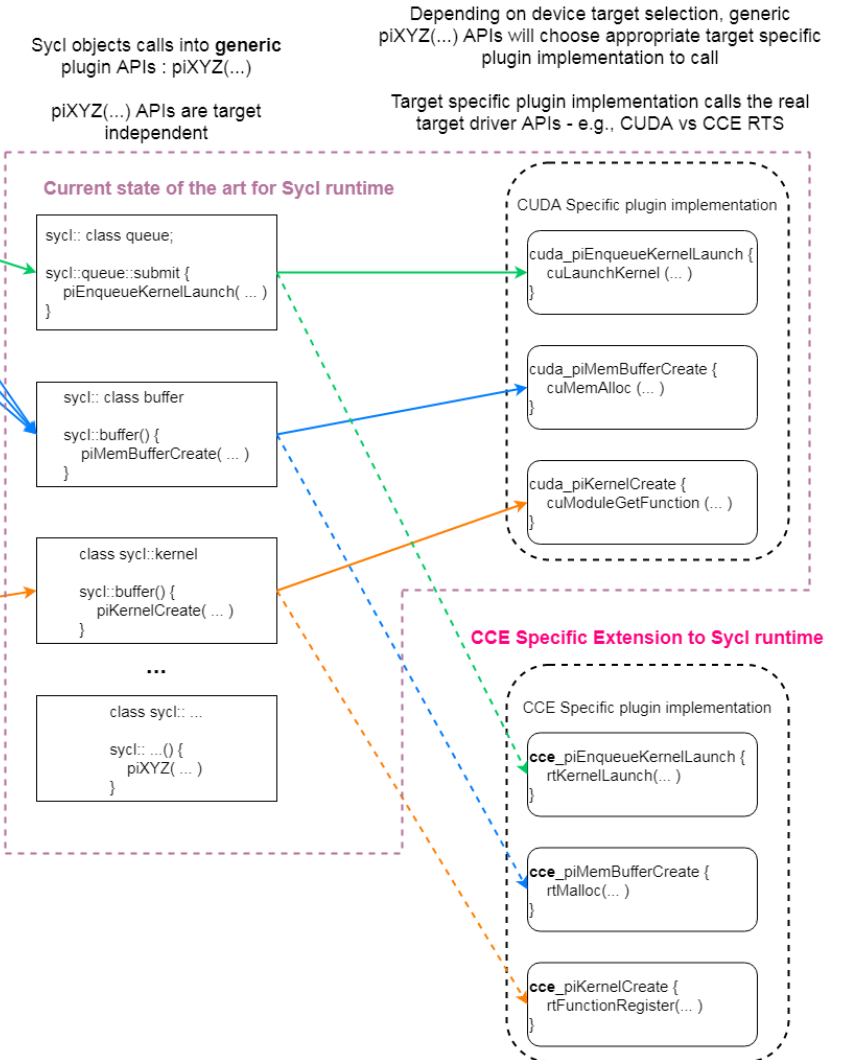
**Current plugins:** OpenCL, Level_zero, CUDA

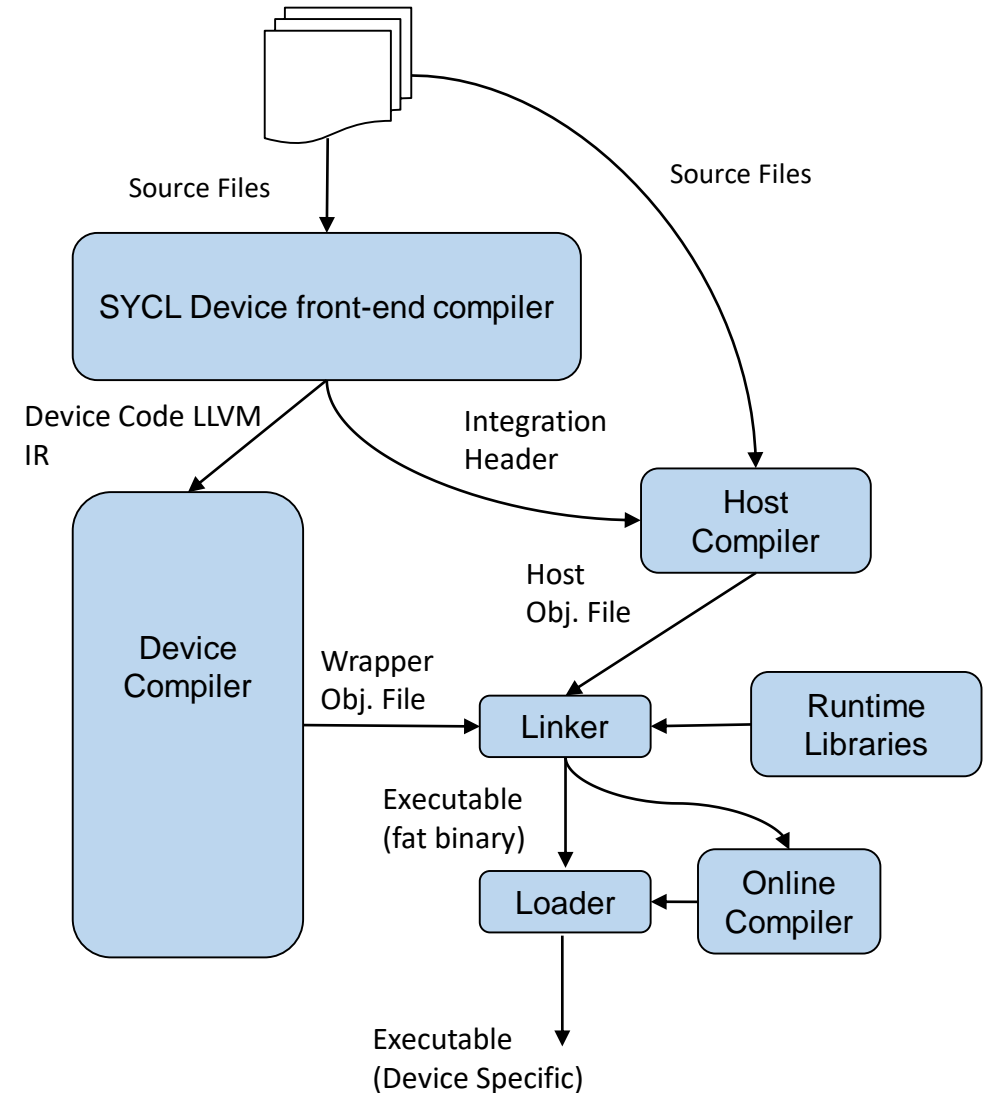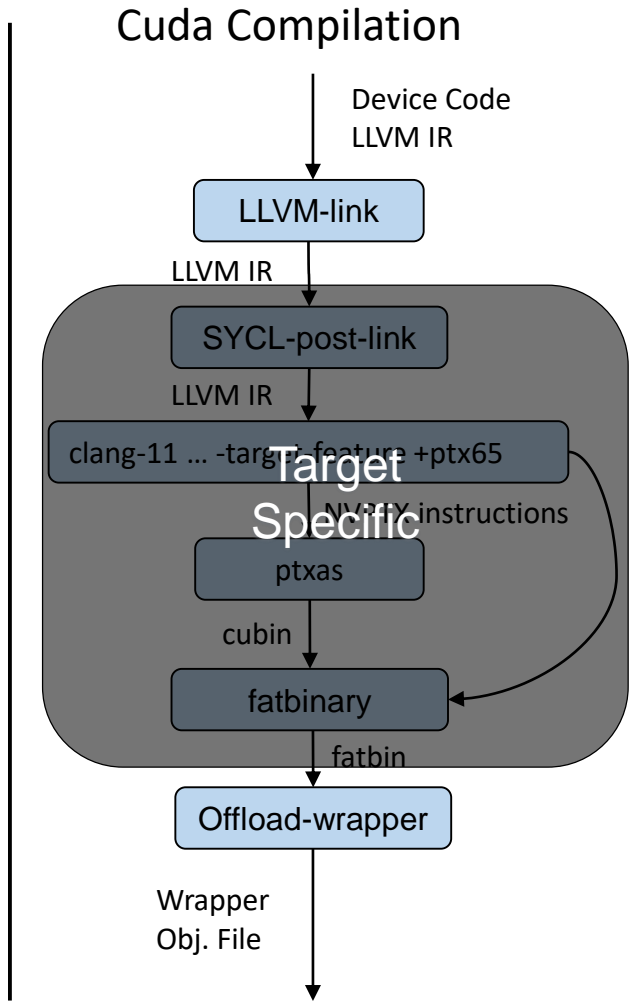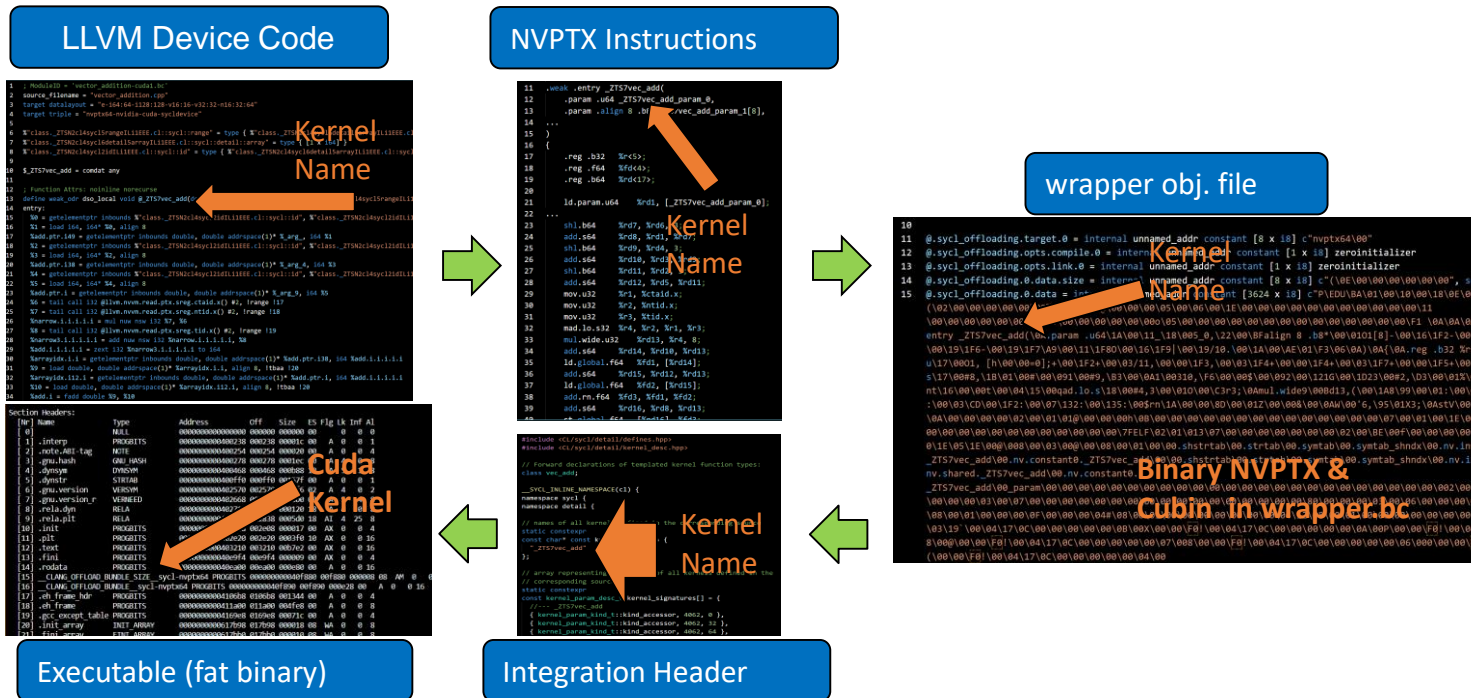**CCE Plugin:** Implementation similar to the CUDA plugin

# Background: Compilation Process in DPC++

- Source codes contain both **host** and **device** code.

- DPC++ compiler **bundles** the **host codes together** and the **device code together**.

- **Device** code is compiled separately and stored in a **wrapper obj. file**.

- **Host** code is compiled with an **Integration Header** containing kernel information

- **Linker** links host and wrapper obj. files + **Runtime Libs.**

- **Loader** loads the **fat binary** for execution. It checks if a target specific executable image exists. If not, the generic image is loaded and **compiled online** to target specific image.

https://github.com/intel/llvm/tree/sycl/sycl/doc

Source Files

Source Files

SYCL Device front-end compiler

Device Code LLVM IR

Integration Header

Host Compiler

Host Obj. File

Device Compiler

Wrapper Obj. File

Linker

Runtime Libraries

Executable (fat binary)

Loader

Online Compiler

Executable (Device Specific)

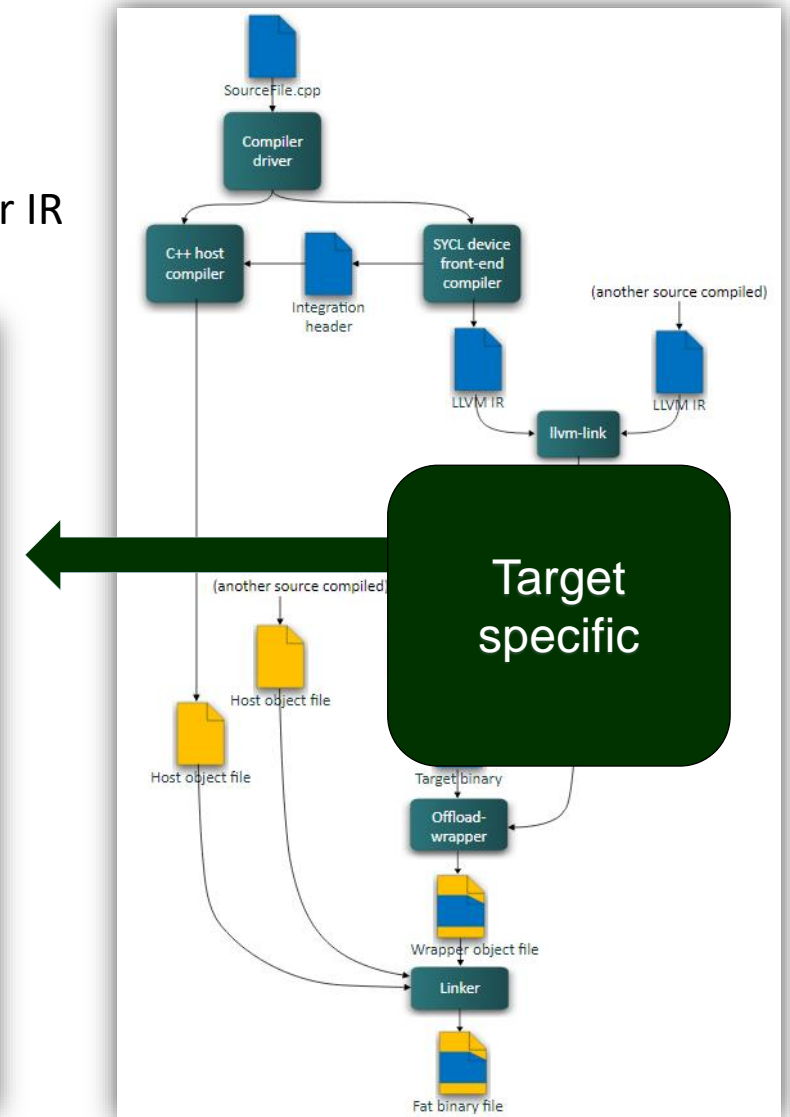# Background: Cuda Device Code Compilation

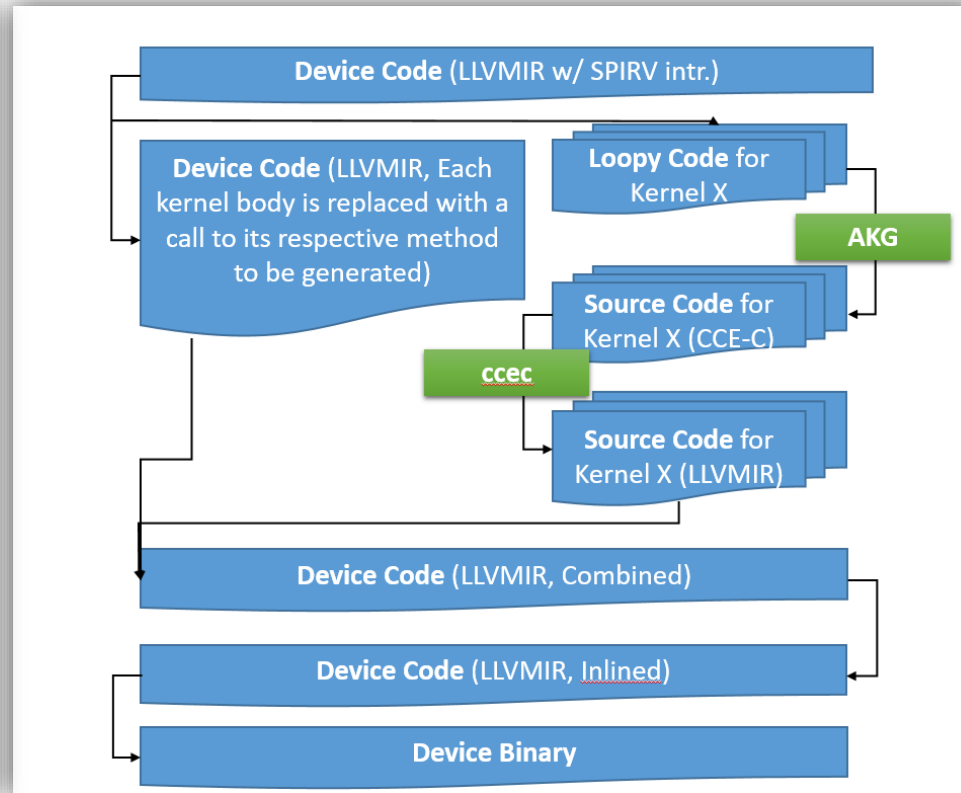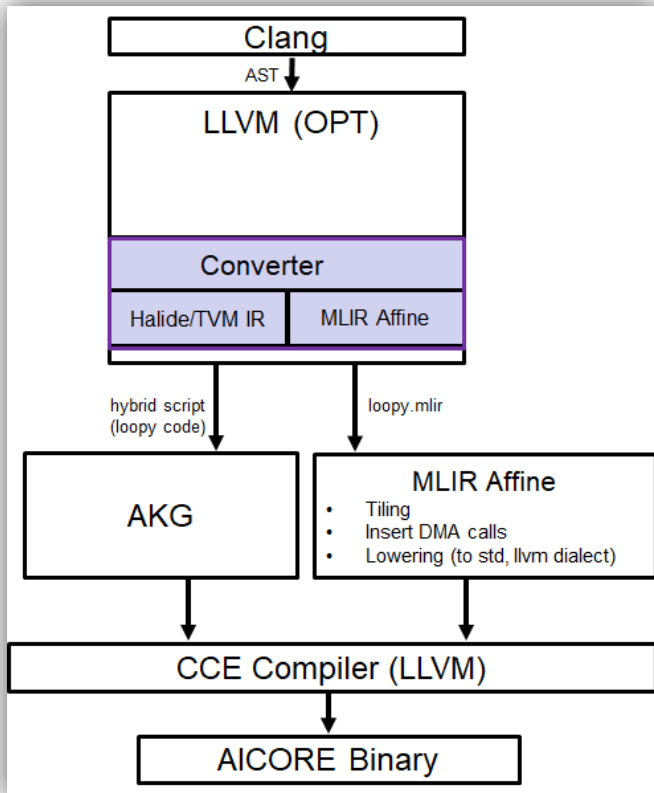- Device Code is compiled to NVPTX and then to Cubin
- This target-specific image is stored in wrapper obj. file
- Online compilation not required

# Our Contribution: CCE Compilation Toolchain

clang++ -fsycl -fsycl-targets=**hiipu64-hisilicon-cce-sycldevice**
simple-hiipu.cpp -o simple-hiipu.exe

- DPC++ device kernel gets compiled into llvm IR and converted to TVM hybrid or IR for AKG. Converter is an llvm opt pass.

# Extension to Unified Shared Memory (USM)

**Implementation of Restricted USM in CCE Plugin:**

1. Map **Host Allocation** to **Host Memory**

2. Map **Shared Memory Allocation** to **Managed Memory (Similar to Cuda)**

3. Map **Device Allocation** to **Device memory**

**PI_CCE Plugin Api to Implement USM:**

- **USM Host Allocation:** piextUSMHostAlloc
- **USM Device Allocation:** piextUSMDeviceAlloc
- **USM Shared Allocation:** piextUSMSharedAlloc
- **USM Free:** piextUSMFree
  - Frees the given USM pointer associated with the context.
- **USM Enqueue Mem. Set:** piextUSMEnqueueMemset
- **USM Enqueue Mem. Copy:** piextUSMEnqueueMemcpy
- **USM Enqueue Mem. Prefetch:** piextUSMEnqueuePrefetch
- **USM Enqueue Mem Advice:** piextUSMEnqueueMemAdvise
  - API to govern behavior of automatic migration mechanisms
- **USM Get Mem. Allocation Info.:** piextUSMGetMemAllocInfo

```
_PI_CL(piEnqueueMemBufferMap, cce_piEnqueueMemBufferMap)
_PI_CL(piEnqueueMemUnmap, cce_piEnqueueMemUnmap)
// USM
_PI_CL(piextUSMHostAlloc, cce_piextUSMHostAlloc)
_PI_CL(piextUSMDeviceAlloc, cce_piextUSMDeviceAlloc)
_PI_CL(piextUSMSharedAlloc, cce_piextUSMSharedAlloc)
_PI_CL(piextUSMFree, cce_piextUSMFree)
_PI_CL(piextUSMEnqueueMemset, cce_piextUSMEnqueueMemset)
_PI_CL(piextUSMEnqueueMemcpy, cce_piextUSMEnqueueMemcpy)
_PI_CL(piextUSMEnqueuePrefetch, cce_piextUSMEnqueuePrefetch)
_PI_CL(piextUSMEnqueueMemAdvise, cce_piextUSMEnqueueMemAdvise)
_PI_CL(piextUSMGetMemAllocInfo, cce_piextUSMGetMemAllocInfo)


_PI_CL(piextKernelSetArgMemObj, cce_piextKernelSetArgMemObj)
_PI_CL(piextKernelSetArgSampler, cce_piextKernelSetArgSampler)
```

# Support for Parallel-for

SIMT to SIMD

we have broken the boundary between the workgroups and workitems:

- AKG (an external tool) is capable of performing **automatic parallelization, vectorization**
- We created our own parallel_for C++ abstraction where we require **static number** of workgroups and workitems
- Converter would artificially create outer loops with an extent that equals to the total number of tasks as well as a loop hint to AKG to select the best parallelization factor
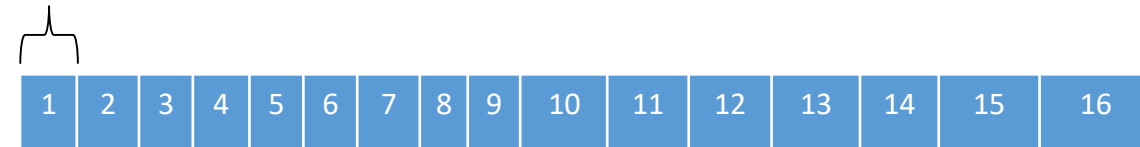
```
struct pi_device_binary_struct {
  uint16_t Version;
  uint8_t Kind;    // 4 for SYCL
  uint8_t Format;  // 1 for native
  const char *DeviceTargetSpec;
  const char *CompileOptions;
  const char *LinkOptions;
  const char *ManifestStart;
  const char *ManifestEnd;
  const unsigned char *BinaryStart;
  const unsigned char *BinaryEnd;
  _pi_offload_entry EntriesBegin;
  _pi_offload_entry EntriesEnd;
  pi_device_binary_property_set PropertySetsBegin;
  pi_device_binary_property_set PropertySetsEnd;
};
```

```
pi_result cce_piextDeviceSelectBinary(pi_device device,
                                      pi_device_binary *binaries,
                                      pi_uint32 num_binaries,
                                      pi_uint32 *selected_binary) {
```
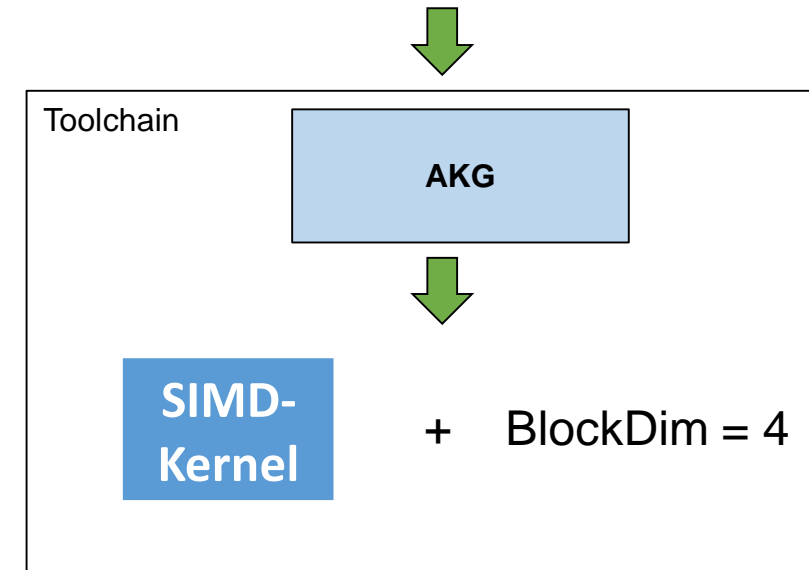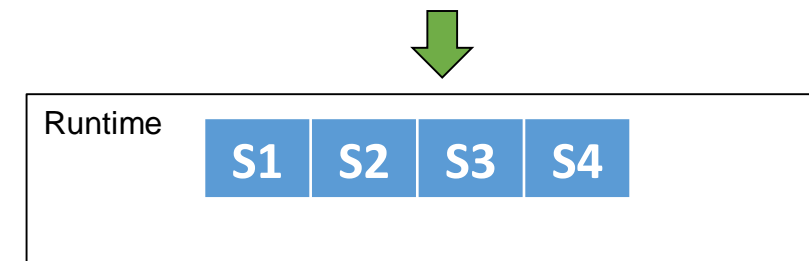
Kernel

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |

- Moving from a SIMT model to a multi-core of SIMD kernels (SIMD-kernels)

Toolchain

**AKG**

**SIMD-Kernel**   +   BlockDim = 4

- At runtime 4 SIMD kernels should be invoked to fully use 4 AICORES

Runtime

| S1 | S2 | S3 | S4 |

# Supported Examples - Matmult

**Size:**

- **M = 384**

- **K = 256**

- **N = 512**

1. Compilation uses C100 sub-architecture of HiIPU target

2. Device selection using hiipu_selector() class

3. Current implementation uses both Vector and Cube units

4. Correctness of output checked on host side

```
$X3CPP_BUILD/bin/clang++ \
    -fsycl \
    -fsycl-targets=hiipu64_c100-hisilicon-cce-sycldevice \
    -DSYCL_DISABLE_PARALLEL_FOR_RANGE_ROUNDING \
    $PROGRAM.cpp \
    -B$AKG_HOME \
    --cce-path=$CCE_HOME
```

```cpp
int main() {
  const size_t M = 6 * 16 * 4;
  const size_t K = 4 * 16 * 4;
  const size_t N = 8 * 16 * 4;

  std::vector<float> A(M * K);
  std::vector<float> B(K * N);
  std::vector<float> C(M * N);

  for (int i = 0; i < M * K; i++) A[i] = 2.0;
  for (int i = 0; i < K * N; i++) B[i] = 2.0;

  queue deviceQueue(hiipu_selector{});
  {

    ndbuffer<float, M, K> A_buf(A.data());
    ndbuffer<float, K, N> B_buf(B.data());
    ndbuffer<float, M, N> OUT_buf(C.data());

    deviceQueue.submit([&](handler &cgh) {
      auto A_acc = A_buf.get_access<sycl_read>(cgh);
      auto B_acc = B_buf.get_access<sycl_read>(cgh);
      auto OUT_acc = OUT_buf.get_access<sycl_write>(cgh);

      auto kern = [=](id<2> ids) {
        size_t m = ids[0];
        size_t n = ids[1];
        float sum = 0.0f;
        for (unsigned int k = 0; k < K; k++) {
          sum += A_acc[m][k] * B_acc[k][n];
        }
        OUT_acc[m][n] = sum;
      };
      static_parallel_for<class matmult, M, N>(cgh, kern);
    });
```

# Future Work

- Alternative codegen support in parallel with AKG
  - MLIR path of code generation is being developed and will require another path to be added into our custom toolchain
- Polish runtime and toolchain code base for production
  - Will be looking forward to upstream our work

**Questions?**