# SYCL 2020 in hipSYCL

## DPC++ features on AMD GPUs, NVIDIA GPUs and CPUs
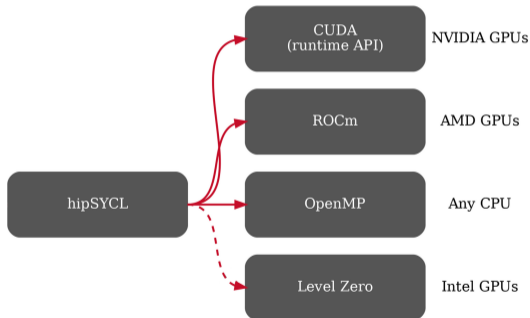
Aksel Alpay

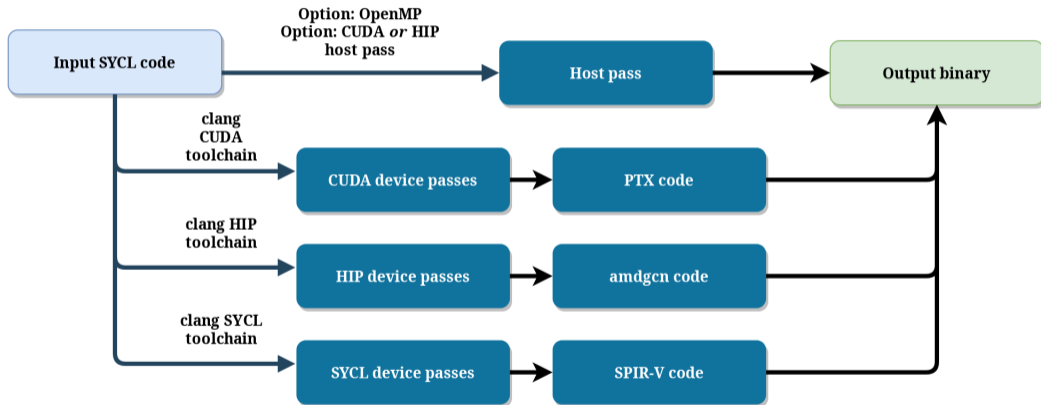Heidelberg University

# Introduction to hipSYCL

**hipSYCL**
A generic, multi-backend SYCL
implementation with emphasis on
aggregating existing toolchains.

► Source-compatible with
   vendor-specific programming
   models
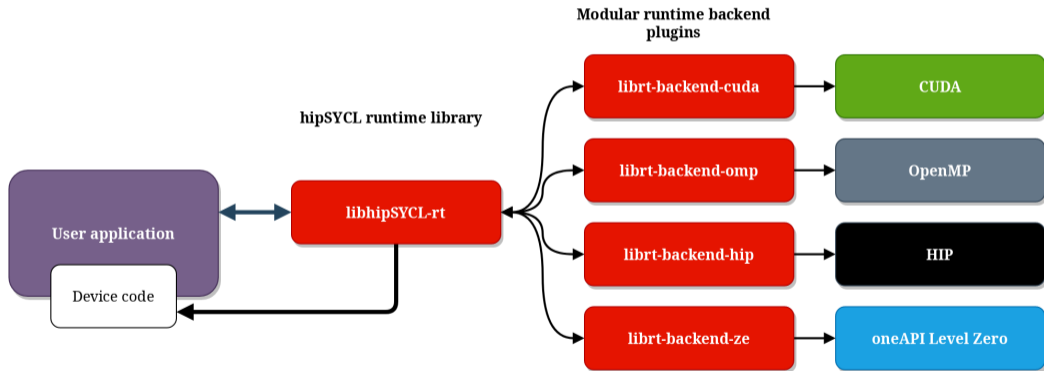► Unique extensions, e.g. full
   buffer-USM interoperability

# hipSYCL: Multiple toolchains in one.

► `syclcc -O3 --hipsycl-targets="omp;cuda:sm_70;hip:gfx906" test.cpp`
► CMake integration available: `find_package(hipSYCL)`, `add_sycl_to_target()`

# hipSYCL runtime architecture

Modular runtime backend plugins

hipSYCL runtime library

User application

Device code

libhipSYCL-rt

librt-backend-cuda → CUDA

librt-backend-omp → OpenMP

librt-backend-hip → HIP

librt-backend-ze → oneAPI Level Zero

# SYCL 2020: Simple things are simple now

```
1  using namespace access;
2  queue q;
3  {
4   buffer<int> a{input_a, size}
5   buffer<int> b{input_b, size};
6   buffer<int> c{output, size};
7   q.submit([&](handler& cgh){
8    auto aa=a.get_access<mode::read>();
9    auto ab=b.get_access<mode::read>();
10   auto ac=c.get_access<mode::write>();
11   cgh.parallel_for<class name>(size
         ,[=](sycl::idx<1> i){
12    ac[i] = aa[i] + ab[i];
13   });
14  });
15 }
```

```
1  queue q;
2  int* a=malloc_shared<int>(size);
3  int* b=malloc_shared<int>(size);
4  int* c=malloc_shared<int>(size);
5  //TODO: Fill input a,b
6  q.parallel_for(size,[=](id<1> i){
7   c[i] = a[i] + b[i];
8  }).wait();
```

# SYCL 2020 in hipSYCL

| | |
|---|---|
| Accessor simplifications | ✔ (partial) (PR) |
| USM: Memory management functions | ✔ (PR) |
| USM: Queue shortcuts | ✔ (PR) |
| USM: Prefetch | ✔ (PR) |
| USM: mem_advise | ✘ |
| USM: memcpy | ✔ (PR) |
| USM: memset/fill | ✔ (PR) |
| host tasks | ✘ |
| Optional lambda naming | ✔ (PR) |
| Subgroups | ✔ (PR) |
| In-order queues | ✔ (PR) |
| Explicit dependencies ( depends_on() ) | ✔ (PR) |
| Backend interop API | ✔ (PR) |
| Reductions | ✔ (PR) |
| Group algorithms | ✔ (PR) |
| New device selector API | ✘ |
| Aspect API | ✘ |
| Deduction guides | ✔ (PR) |
| atomic_ref | ✘ |
| marray | ✘ |
| New SYCL/sycl.hpp header | ✔ (PR) |
| C++17 by default | ✔ (PR) |

| | |
|---|---|
| Builtin changes: ctz() , clz() | ✘ |
| Remove *_class types | ✘ |
| const return type for read accessor operator[] | ✘ |
| Remove buffer API for unique_ptr | ✘ |
| Replace program class with module | ✘ |
| Add kernel_handler | ✘ |
| explicit queue , context constructors | ✔ (PR) |
| Only require C++ trivially copyable for shared data | ✔ |
| Update group class with new types/member functions | ✘ |
| Remove nd_item::barrier() | ✘ |
| Replace mem_fence with atomic_fence | ✘ |
| Add vec::operator[] ,unary +,- . static constexpr get_size()/get_count() | ✔ (PR) |
| buffer, local accessor are C++ ContiguousContainer | ✘ |
| Replace image with sampled_image , unsampled_image | ✘ |
| All accessors are placeholders | ✔ (PR) |
| Use single exception type derived from std::exception | ✘ |
| Default asynchronous handler should terminate program | ✔ (PR) |

| | |
|---|---|
| Kernel invocation APIs take const reference to kernels, kernels must be immutable | ✘ |
| Queue constructor accepting both device and context | ✘ |
| Simplified parallel_for API | ✘ |
| Clarified names for device specific info queries | ✘ |
| Address space changes, generic address spaces | ✘ |
| Updated multi_ptr interface | ✘ |
| Remove OpenCL types, cl_int etc | ✔ |

https://github.com/hipSYCL/featuresupport

# Key DPC++/SYCL 2020 features implemented

Our work as oneAPI Center of Excellence:

- ▶ Unified Shared Memory (USM)
- ▶ Optional Lambda Naming
- ▶ Subgroups
- ▶ Parallel algorithms for groups and subgroups
- ▶ Parallel reductions
- ▶ Queue shortcuts
- ▶ Explicit task graphs
- ▶ …

**hipSYCL support increases adoption and portability of SYCL 2020 features (e.g. AMD GPUs)**
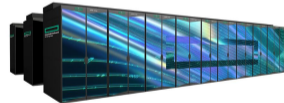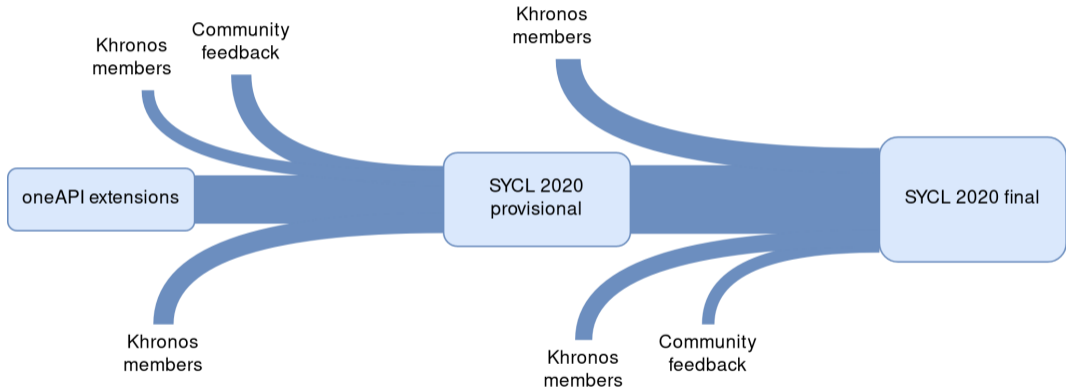
# Frontier, El Capitan, LUMI

**Bring SYCL 2020 to upcoming supercomputers based on AMD GPUs**

# From DPC++ extensions to SYCL 2020



- ▶ SYCL 2020 interfaces in current implementations and client code may vary
- ▶ hipSYCL interfaces started with SYCL 2020 provisional, now moving towards SYCL 2020 final

# Unified Shared Memory in hipSYCL

- ▶ Device-accessible host memory
- ▶ Explicit USM
- ▶ Shared allocations (on-demand migration)
- ▶ Performance hints (prefetch/memadvise)

```
1  sycl::queue q;
2  int* ptr =
3   sycl::malloc_shared<int>(size,q)
4  q.parallel_for(size,
5    [=](sycl::id<1> idx){
6    const int i = idx.get(0);
7    ptr[i] = i;
8  });
```

| CUDA | HIP | CPU |
|---|---|---|
| cudaMallocHost | hipHostMalloc | (regular host allocations) |
| cudaMalloc | hipMalloc | |
| cudaMallocManaged | hipMallocManaged | |
| cudaMemPrefetchAsync | hipMemPrefetchAsync | |
| cudaMemAdvise | hipMemAdvise | |

# hipSYCL USM performance

## Parallel research kernels benchmarks[1]



ROCm does not yet fully support allocations with on-demand page migration.

[1] https://github.com/ParRes/Kernels

# Subgroups

- ▶ Expose hardware below work group granularity
- ▶ SIMD units
- ▶ Useful for optimization

```
1  sycl::nd_item<1> idx = ...;
2  auto sgrp = idx.get_sub_group();
```

**CUDA**
Mapped to CUDA warps

**HIP**
Mapped to AMD wavefronts

**CPU**
Individual subgroup for each work item

- ▶ Difficult for library-only CPU backends
- ▶ Vectorization controlled by OpenMP compiler

# Group algorithms in hipSYCL

- `any_of`, `none_of`, `all_of`
- Reductions, scans
- broadcast, barrier
- At work group and subgroup level
- Collective and Iterator-based variants

```
1  int myval = ...;
2  int sum=sycl::reduce_over_group(
3    group, myval, sycl::plus<int>{})
```

| CUDA | HIP | CPU |
|---|---|---|
| ▶ Subgroup intrinsics, warp shuffles<br>▶ Optimized local memory usage | ▶ Subgroup intrinsics, warp shuffles<br>▶ Optimized local memory usage | ▶ Sequential with OpenMP vectorization<br>▶ Bound by synchronization |

# Group reduce and scan on AMD Radeon VII

Work group reduce



Work group inclusive scan

► Competitive performance compared to rocPRIM

► We are handicapped: Group size not known at compile time

Wünsche, H. (2021): "SYCL 2020 work group parallel primitives [...] in hipSYCL". Bachelor Thesis. Heidelberg U.

# Parallel reductions

- ▶ Variadic scalar reductions
- ▶ basic, `nd_range`, hierarchical and scoped parallelism supported
- ▶ No multi-dimensional reductions yet
- ▶ Huge optimization space!

```
q.parallel_for(range{size},
 reduction(output,
           sycl::plus<int>{}),
 [=](sycl::id<1> idx,
     auto& reducer){
  int myvalue = ...;
  reducer += myvalue;
});
```

**CUDA**

- ▶ Work group reductions
- ▶ Multiple kernel launches

**HIP**

- ▶ Work group reductions
- ▶ Multiple kernel launches

**CPU**

- ▶ Per-thread reductions to cache-line aligned private storage
- ▶ Cross-thread final reduction

# Reduction performance



Reducing $10^9$ ints on AMD Radeon VII



Reducing $10^9$ ints on Intel Xeon Gold 6130

▶ hipSYCL reduction performance can compete with vendor-optimized libraries

▶ SYCL model is very flexible and allows for more user control – and user error!

▶ **For pure reductions, we still need optimized SYCL libraries!**

# oneAPI libraries

**The language is only half the way for oneAPI portability! We need the libraries.**

▶ oneMKL: hipSYCL (and rocBLAS) support currently WIP
▶ oneDPL: Most unit tests already run with hipSYCL[2]



[2]https://github.com/hipSYCL/oneDPL

# Case study: oneDPL

**What was necessary to port oneDPL to hipSYCL?**

- ▶ Build system
- ▶ SYCL compiler detection/macros
- ▶ pre-SYCL 2020 final APIs:
    - ▶ `ONEAPI` namespace for reductions/group algorithms
    - ▶ `noinit` vs `no_init` property
- ▶ DPC++ implementation details
    - ▶ `mode_tag_t`
- ▶ Optimize work group size selection for hipSYCL

- ▶ **No functional changes**
- ▶ **Single file with some compatibility aliases**

# Conclusion

UNIVERSITÄTS-
RECHENZENTRUM

UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

# hipSYCL

▶ Rapidly moving towards SYCL 2020
▶ Key features supported - high performance implementations for all backends.

**It _is_ possible to write standard SYCL 2020/DPC++ code, be portable without sacrificing performance w.r.t vendor-optimized libraries!**

▶ Make SYCL 2020 ubiquitous - let DPC++/SYCL 2020 code run on AMD GPUs, NVIDIA GPUs, any CPU
▶ Ongoing work: oneMKL, oneDPL, other SYCL 2020 features, more optimizations…
▶ Open source: `https://github.com/illuhad/hipSYCL`
▶ Get in touch: `aksel.alpay@uni-heidelberg.de`