# Development and optimization of a SYCL backend for libCEED

## Umesh Unnikrishnan[1], Kris Rowe[1], and Varsha Madananth[2]

[1] Argonne Leadership Computing Facility
[2] Intel Corporation

# Outline

➢ **Software design of libCEED**

➢ **SYCL online compiler**

➢ **Optimization of hotspot kernels**

  ❑ Specialization constant

  ❑ Tuning workgroup sizes/barriers

  ❑ SIMD size/register width

➢ **Summary and Future work**

**Disclaimer: This work was done on a pre-production supercomputer with early versions of the Aurora software development kit.**

Argonne
NATIONAL LABORATORY

# libCEED

Developed through CEED co-design center as part of the ECP courtesy: Jed Brown, Natalie Beams and others.

https://github.com/CEED/libCEED

## Library

- C language
- Element-based discretizations
- low & high-order

## Interfaces

- Fortran
- Python
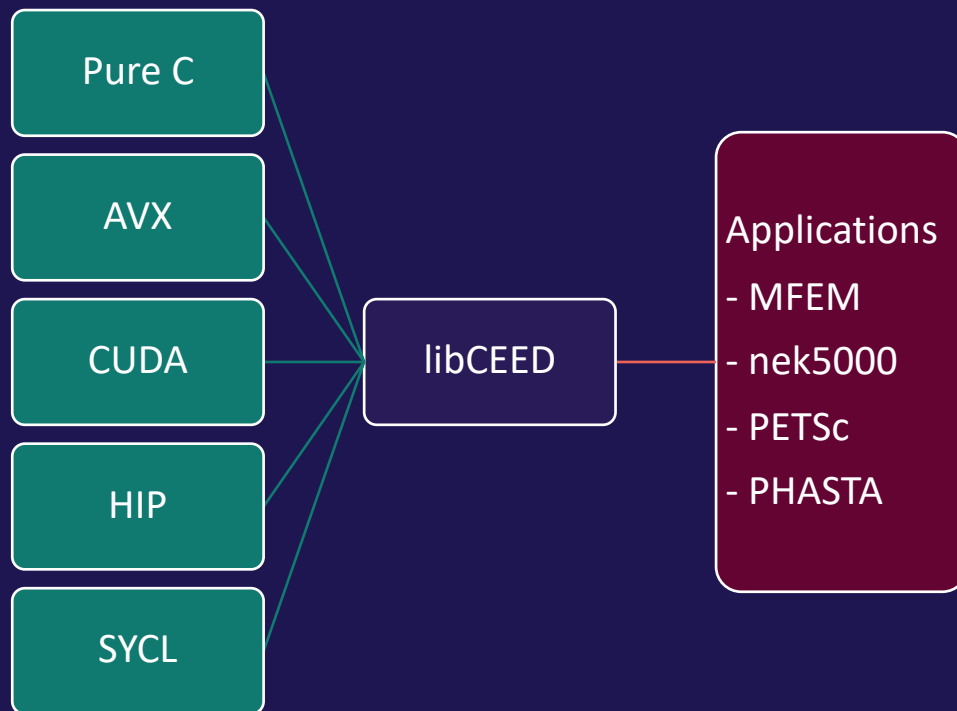- Julia
- Rust

## Backends

- Runtime selection
- CPU—serial, AVX, LIBXSMM
- GPU (Native)—CUDA, HIP
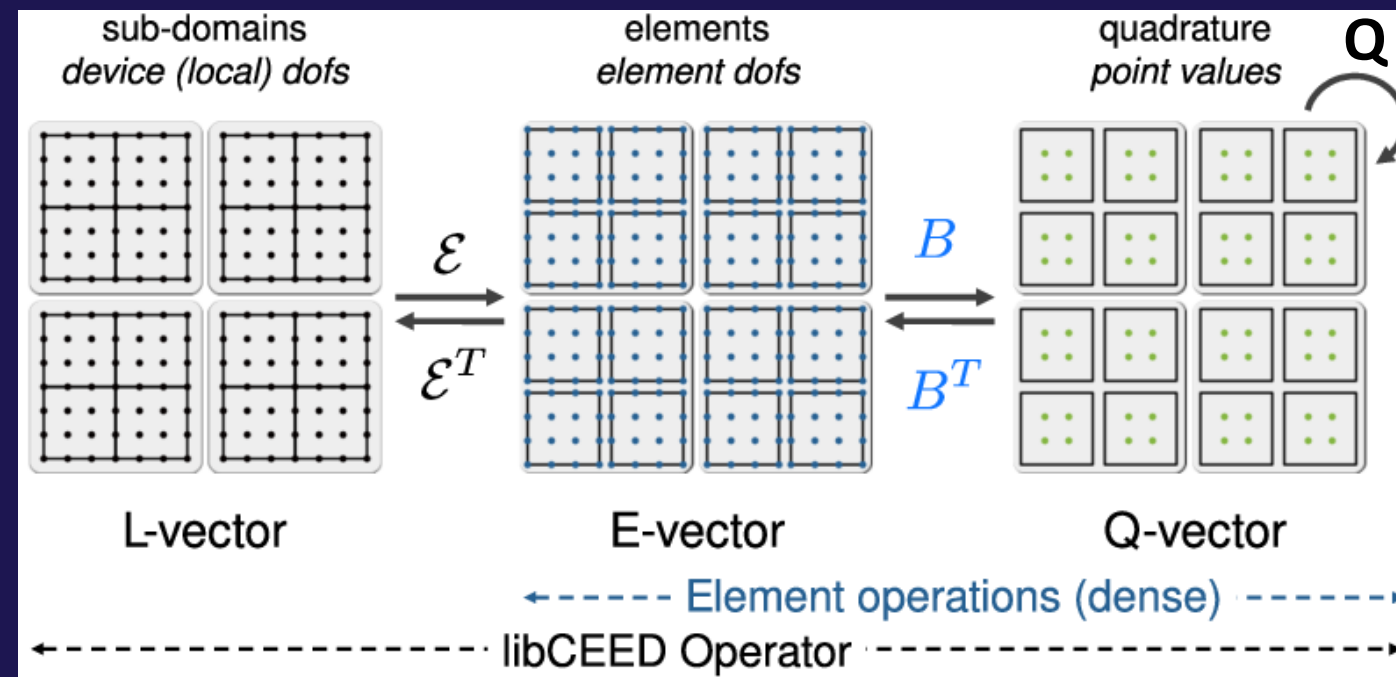- MAGMA, OCCA**

## Implementation

- PIMPL Idiom
- Host code is serial
- User handles MPI

Argonne
NATIONAL LABORATORY

# Overview of libCEED operations

- Portable library that provides an API for applications to share efficient kernels for element-based discretizations.

$$A = E^T B^T Q B E$$

*https://github.com/CEED/libCEED*

# libCEED Runtime Compilation Usage

| Problem Parameters | Constants | polynomial order |
| --- | --- | --- |
| | | spatial dimension |
| | Determine loop-bounds | element count |
| | | node count |
| User Functions | Operate pointwise at mesh nodes | Boundary conditions |
| | | Forces |
| | Injected into kernels | Problem-dependent |
| | | Avoids memory traffic |

Argonne
NATIONAL LABORATORY

# Runtime code compilation in libCEED

- libCEED host application generates device code at runtime based on user input (the code is in stringstream) – compiles and runs on GPU devices.

- CUDA backend uses nvrtc for runtime compilation of generated kernels.

- Two possible solutions for SYCL –

  ➢ Using specialization constant (a partial alternative)

  ➢ Using Intel's oneAPI online compilation extension (experimental early access) –

    https://github.com/intel/llvm/blob/sycl/sycl/doc/extensions/experimental/sycl_ext_intel_online_compiler.asciidoc

Argonne
NATIONAL LABORATORY

# Existing CUDA Backend

AOT Compilation

- Uses nvcc
- Some (few) kernels

Online Compilation

- Uses NVRTC
- Kernel source string is generated
- Boilerplate code is provided (e.g., read/writes with global mem)
- User source code is loaded from file
- Create  BP→UserDefinedFunc→BP "sandwich"
- Load Cumodule, use CUfunction

Argonne
NATIONAL LABORATORY

# Designing the SYCL Backend

| Problem Parameters | Specialization constants |
| --- | --- |
| | Define kernels using lambdas or functors |
| | CUDA module can be replaced by SYCL kernel bundle |
| User Functions | How to compile source code string? (Ideal) |
| | How to link compiled user functions?  (Fallback) |

Argonne
NATIONAL LABORATORY

# libCEED SYCL Online Compilation

- The libCEED CUDA and HIP backends use their respective vendor runtime compilation libraries (e.g., nvrtc, hiprtc).

- The current SYCL spec doesn't prescribe similar functionality.

- An Intel extension for SYCL allows for the runtime compilation of OpenCL C source.

  - Provided to SYCL API as a std::string

- Restriction to OpenCL C required some workarounds when porting libCEED jit source.

  - E.g., No templates, no function pointers

  - Also required some workarounds in QFunction implementation.

- We are working with Intel compiler team to help drive the online compiler extension forward.

- Future versions may support runtime compilation of SYCL source

  https://github.com/intel/llvm/blob/sycl/sycl/doc/extensions/experimental/sycl_ext_oneapi_kernel_compiler.asciidoc

Argonne
NATIONAL LABORATORY

# Intel Online Compiler SYCL Extension

- Experimental extension

- Calls ocloc library API (libigc)

- Compile OpenCL string to "binary blob"

- How to use "binary blob"?

Alternatives are not ideal

- OpenCL plugin only + interop

- Use Level Zero + libigc directly

```cpp
#include "sycl/ext/intel/online_compiler.hpp"

#include <iostream>
#include <vector>

static const char *kernelSource = R"===(
__kernel void my_kernel(__global int *in, __global int *out) {
  size_t i = get_global_id(0);
  out[i] = in[i] + 1;
}
)===";

using namespace sycl::INTEL;

int main(int argc, char **argv) {
  online_compiler<source_language::opencl_c> compiler;
  std::vector<byte> blob;

  try {
    blob = compiler.compile(
      std::string(kernelSource),
      std::vector<std::string> {
        std::string("-cl-fast-relaxed-math")
      }
    );
  }
  catch (online_compile_error &e) {
    std::cout << "compilation failed\n";
    return 1;
  }
  return 0;
}
```

Example from Intel LLVM docs

Argonne
NATIONAL LABORATORY

# Example

Online Compiler + Level Zero

```cpp
sycl_queue.submit([&](sycl::handler &cgh) {
  cgh.depends_on({copy_x, copy_y});

  cgh.set_args(alpha, x, y);
  cgh.parallel_for(N, *axpy_kernel);
});
```

SYCL kernel bundle

SYCL kernel

```cpp
sycl::kernel *buildKernelFromSource(const sycl::device &sycl_device,
                                    const sycl::context &sycl_context,
                                    const std::string &kernel_name,
                                    const std::string &kernel_source,
                                    const std::vector<std::string> &flags) {
  using namespace sycl::ext::intel::experimental;
  online_compiler<source_language::opencl_c> compiler(sycl_device);

  std::vector<unsigned char> blob;
  try { blob = compiler.compile(kernel_source, flags);}
  catch (online_compile_error &e) {
    std::cout << "compilation failed\n";
    std::exit(EXIT_FAILURE);
  }
  auto lz_device = sycl::get_native<sycl::backend::ext_oneapi_level_zero>(sycl_device);
  auto lz_context = sycl::get_native<sycl::backend::ext_oneapi_level_zero>(sycl_context);

  const char *lz_flags = "-ze-opt-level=2";
  ze_module_desc_t lz_mod_desc = {ZE_STRUCTURE_TYPE_MODULE_DESC,
                                  nullptr,
                                  ZE_MODULE_FORMAT_IL_SPIRV,
                                  blob.size(),
                                  blob.data(),
                                  lz_flags,
                                  nullptr};
  ze_module_handle_t lz_module;
  zeModuleCreate(lz_context, lz_device, &lz_mod_desc, &lz_module, nullptr);
  auto *sycl_module = new sycl::kernel_bundle<sycl::bundle_state::executable>(
    sycl::make_kernel_bundle<sycl::backend::ext_oneapi_level_zero,sycl::bundle_state::executable>(
      {lz_module, sycl::ext::oneapi::level_zero::ownership::transfer},sycl_context
  ));

  ze_kernel_desc_t lz_kernel_desc = {ZE_STRUCTURE_TYPE_KERNEL_DESC, nullptr, 0,kernel_name.c_str()};
  ze_kernel_handle_t lz_kernel;
  zeKernelCreate(lz_module, &lz_kernel_desc, &lz_kernel);
  auto *sycl_kernel = new sycl::kernel(
    sycl::make_kernel<sycl::backend::ext_oneapi_level_zero>(
      {*sycl_module, lz_kernel,sycl::ext::oneapi::level_zero::ownership::transfer},sycl_context
  ));

  return sycl_kernel;
}
```

binary blob

LZ module descriptor

LZ module

LZ kernel

Argonne NATIONAL LABORATORY

# libCEED Implementation

```cpp
using SyclModule_t = sycl::kernel_bundle<sycl::bundle_state::executable>;

int CeedJitGetKernel_Sycl(Ceed ceed, const SyclModule_t *sycl_module, const std::string &kernel_name, sycl::kernel **sycl_kernel) {
  Ceed_Sycl *data;
  CeedCallBackend(CeedGetData(ceed, &data));

  // sycl::get_native returns std::vector<ze_module_handle_t> for lz backend
  // https://github.com/intel/llvm/blob/sycl/sycl/doc/extensions/supported/sycl_ext_oneapi_backend_level_zero.md
  ze_module_handle_t lz_module = sycl::get_native<sycl::backend::ext_oneapi_level_zero>(*sycl_module).front();

  ze_kernel_desc_t   lz_kernel_desc = {ZE_STRUCTURE_TYPE_KERNEL_DESC, nullptr, 0, kernel_name.c_str()};
  ze_kernel_handle_t lz_kernel;
  zeKernelCreate(lz_module, &lz_kernel_desc, &lz_kernel);

  *sycl_kernel = new sycl::kernel(sycl::make_kernel<sycl::backend::ext_oneapi_level_zero>(
      {*sycl_module, lz_kernel, sycl::ext::oneapi::level_zero::ownership::transfer}, data->sycl_context));

  return CEED_ERROR_SUCCESS;
}
```

Argonne
NATIONAL LABORATORY

# SYCL queue synchronization for external libraries

- Streams in CUDA and HIP are in-order.

- A global default stream is available - simplifying the coordination with external libraries.

- Initial implementation of SYCL backends used in-order queues.

- CEED-PHASTA uses PETSc

    - For Intel GPUs, PETSc uses Kokkos-SYCL backend

    - Kokkos-SYCL uses out-of-order SYCL queues

- For synchronization with PETSc, libCEED SYCL inputs SYCL queue from PETSc.

- To use out-of-order queues with the libCEED SYCL implementation, the Intel SYCL extension for enqueueing (asynchronous, device-side) enqueue barriers were used.

- Future development will look to use SYCL events to explicitly express data dependencies within libCEED.

Argonne
NATIONAL LABORATORY

# Performance Baseline

- sycl-fluids examples using Blasius test case.

    - Overall time per iteration on A100 was 3.7 sec

    - Overall time per iteration on PVC was 70 seconds

- 3 hotspots identified for both A100 and PVC.

    - CeedBasisSyclGrad – SYCL kernel

    - CeedBasisSyclInterp – SYCL kernel

    - IJacobian (OpenCL kernel which is online compiled)

| Kernel | PVC (1 Tile) | A100 (CUDA) |
|--------|--------------|-------------|
| IJacobian | 29 ms | 0.92 ms |
| Interp | 10.7 ms | 0.11 ms |
| Grad | 197 ms | 1.3 ms |

Argonne
NATIONAL LABORATORY

# BASIS GRAD KERNEL

Using microbenchmarks of the kernel

**Inline kernel**

```
7    #define CeedScalar float
8    #define NTRIALS 1001
9
10   class CeedBasisSyclGrad;
11
12   int main() {
13     // Setup options
14     const CeedInt Q_1d = 2;
15     const CeedInt num_elem = 10*190*141;
16     const CeedInt dim = 3;
17
18     // Device memory allocations and host to device copy
19
20     // Kernel setup and execution
21     sycl::range<1>      local_range(work_group_size);
22     sycl::range<1>      global_range(num_elem * work_group_size);
23     sycl::nd_range<1>   kernel_range(global_range, local_range);
24
25     float walltimes[NTRIALS], avg_time;
26     for(int trial=0;trial<NTRIALS;trial++) {
27       auto start_time = std::chrono::high_resolution_clock::now();
28
29       sycl_queue.submit([&](sycl::handler &cgh) {
30         sycl::local_accessor<CeedScalar> s_mem(2 * (P * Q + buf_len), cgh);
31
32         cgh.parallel_for<CeedBasisSyclGrad>(kernel_range, [=](sycl::nd_item<1> work_ite
33           .
34           .
35         });
36       });
37       sycl_queue.wait_and_throw();
38
39       auto finish_time = std::chrono::high_resolution_clock::now();
40       walltimes[trial] = std::chrono::duration<float,std::milli>(finish_time-start_time
41   }
```

**Kernel execution time = 1.33 ms**

**Kernel in function**

```
26   int CeedBasisGrad_Sycl(sycl::queue &sycl_queue, CeedBasis_Sycl *impl, CeedScalar *u, CeedScalar *v) {
27     const CeedInt dim          = impl->dim;
28     const CeedInt Q_1d         = impl->Q_1d;
29     const CeedScalar *grad_1d  = impl->d_grad_1d;
30
31     const CeedInt work_group_size    = 32;
32     sycl::range<1>      local_range(work_group_size);
33     sycl::range<1>      global_range(num_elem * work_group_size);
34     sycl::nd_range<1>   kernel_range(global_range, local_range);
35
36     sycl_queue.submit([&](sycl::handler &cgh) {
37       sycl::local_accessor<CeedScalar> s_mem(2 * (P * Q + buf_len), cgh);
38       cgh.parallel_for<CeedBasisSyclGrad>(kernel_range, [=](sycl::nd_item<1> work_item) {
39         ...
40       });
41     });
42     return 1;
43   }
44
45   int main() {
46     // Setup options
47     const CeedInt Q_1d = 2;
48     const CeedInt dim = 3;
49
50     // SYCL initialization
51
52     // Ceed Basis setup
53     CeedBasis_Sycl *basis = (CeedBasis_Sycl*)calloc(1,sizeof(CeedBasis_Sycl));
54     basis->dim = dim;
55     basis->Q_1d = Q_1d;
56     basis->d_grad_1d = sycl::malloc_device<CeedScalar>(interp_length, sycl_device, sycl_context);
57
58     for(int trial=0;trial<NTRIALS;trial++) {
59       auto start_time = std::chrono::high_resolution_clock::now();
60       result = CeedBasisGrad_Sycl(sycl_queue, basis, d_u, d_v);
61       sycl_queue.wait_and_throw();
62
63       auto finish_time = std::chrono::high_resolution_clock::now();
64       walltimes[trial] = std::chrono::duration<float,std::milli>(finish_time-start_time).count();
65   }
```

**Kernel execution time = 5.8 ms**

Argonne
NATIONAL LABORATORY

# Use of SYCL specialization constants

```cpp
12  using SpecID = sycl::specialization_id<CeedInt>;
13  using SyclModule_t = sycl::kernel_bundle<sycl::bundle_state::executable>;
14
15  static constexpr SpecID BASIS_Q_1D_ID;
16  static constexpr SpecID BASIS_DIM_ID;
17
18  int CeedBasisGrad_Sycl (sycl::queue &sycl_queue, const SyclModule_t& sycl_module, CeedBasis_Sycl *impl, const CeedScalar *u, CeedScalar *v) {
19      sycl::range<1>        local_range(work_group_size);
20      sycl::range<1>        global_range(num_elem * work_group_size);
21      sycl::nd_range<1>    kernel_range(global_range, local_range);
22
23      sycl::event e = sycl_queue.ext_oneapi_submit_barrier();
24      sycl_queue.submit([&](sycl::handler &cgh) {
25          cgh.depends_on({e});
26          cgh.use_kernel_bundle(sycl_module);
27          sycl::local_accessor<CeedScalar> s_mem(2 * (op_len + buf_len), cgh);
28
29          cgh.parallel_for<CeedBasisSyclGrad<transpose>>(kernel_range, [=](sycl::nd_item<1> work_item, sycl::kernel_handler kh) {
30              // Retrieve spec constants ------------------------------------->
31              const CeedInt dim           = kh.get_specialization_constant<BASIS_DIM_ID>();
32              const CeedInt Q_1d          = kh.get_specialization_constant<BASIS_Q_1D_ID>();
33
34              // Work
35          });
36      });
37      return 1;
38  }
39
40  int main() {
41      // Initialization
42
43      std::vector<sycl::kernel_id> kernel_ids = { sycl::get_kernel_id<CeedBasisSyclGrad<1>>(), sycl::get_kernel_id<CeedBasisSyclGrad<0>>() };
44      sycl::kernel_bundle<sycl::bundle_state::input> input_bundle = sycl::get_kernel_bundle<sycl::bundle_state::input>(sycl_context, kernel_ids);
45
46      input_bundle.set_specialization_constant<BASIS_DIM_ID>(dim);
47      input_bundle.set_specialization_constant<BASIS_Q_1D_ID>(Q_1d);
48      SyclModule_t *sycl_module = new SyclModule_t(sycl::build(input_bundle));
49
50      for(int trial=0;trial<NTRIALS;trial++) {
51          auto start_time = std::chrono::high_resolution_clock::now();
52          result = CeedBasisGrad_Sycl<1>(sycl_queue, *sycl_module, num_elem, basis, d_interp_1d, d_grad_1d, d_u, d_v);
53          sycl_queue.wait_and_throw();
54
55          auto finish_time = std::chrono::high_resolution_clock::now();
56          walltimes[trial] = std::chrono::duration<float,std::milli>(finish_time-start_
57      }
```

Kernel execution time = 1.38 ms

**Argonne** NATIONAL LABORATORY

# Further optimizations on CeedBasisGrad

- Adjust WG size from 1024 to 32.

- Replace workgroup barrier with nd_item barrier

  sycl::group_barrier(work_group) → work_item.barrier(sycl::access::fence_space::local_space)

- Remove IGC runtime memory checks by using specialization constant.

  warning: from dir:/home/vmadananth/PHASTA_aesp_CNDA/CEED-PHASTA/libCEED/backends/sycl-ref from file:ceed-sycl-ref-basis.sycl.cpp line:100 :Adding additional control flow due to presence of generic address space operations
  warning: from dir:/home/vmadananth/PHASTA_aesp_CNDA/CEED-PHASTA/libCEED/backends/sycl-ref from file:ceed-sycl-ref-basis.sycl.cpp line:185 :Adding additional control flow due to presence of generic address space operations
  warning: from dir:/home/vmadananth/PHASTA_aesp_CNDA/CEED-PHASTA/libCEED/backends/sycl-ref from file:ceed-sycl-ref-basis.sycl.cpp line:187 :Adding additional control flow due to presence of generic address space operations
  warning: from dir:/home/vmadananth/PHASTA_aesp_CNDA/CEED-PHASTA/libCEED/backends/sycl-ref from file:ceed-sycl-ref-basis.sycl.cpp line:0 :Adding additional control flow due to presence of generic address space operations

| Kernel | PVC (1 Tile) | A100 |
|---|---|---|
| IJacobian | 29ms | 0.9ms |
| Interp | 1.75ms | 0.11ms |
| Grad | 4.8ms | 1.3ms |

Argonne
NATIONAL LABORATORY

# Optimizations performed on IJacobian

- IJacobian kernel was using default range. The compiler was setting the workgroup size to 32.
  - Changing to default nd_range and setting WG size to 384 (CUDA) improved performance slightly.

- Reducing register spills
  - Building with AOT did not show any warnings on spills as this was OpenCL kernel.
  - Inspected assembly by setting the following env variables –
    - IGC_DumpToCurrentDir=1,  IGC_DumpToCurrentDir=1
    - For more IGC env variables - https://github.com/intel/intel-graphics-compiler/blob/master/documentation/configuration_flags.md
    - Generates asm for all kernels. To search for kernel names –
      
      for f in ./*.asm; do echo "------------"; echo $f; cat $f | grep "\/\/.kernel"; done

```
//.kernel CeedKernelSyclRefQFunction_IJacobian_Newtonian_Prim
//.platform PVCXT
//.thread_config numGRF=128, numAcc=4, numSWSB=16
//.options_string ""
//.full_options "-emitLocation -forceAssignRhysicalReg "" -hasRNEandDenorm -noStitchExternFunc -linker 0 -lscEnableImmOffsFor 196638 -
preserver0 -TotalGRFNum 128 -abortOnSpill 4 -boundsChecking -presched-ctrl 6 -presched-rp 100 -nodpsendreorder -SBIDDepLoc -output -
binary -dumpcommonisa -shaderDumpFilter "" -dumpvisa -printHexFloatInAsm -noverifyCISA -enableHalfLSC -hasInt64Add -partialInt64 -
generateDebugInfo "
//.instCount 3207
//.RA type  GRAPH_COLORING_SPILL_FF_RA
//.spill size 54784
//.spill GRF est. ref count 1475
```

Argonne NATIONAL LABORATORY

# Reducing spills

- By default , PVC has 128 64-byte registers allocated per thread.

- Registers spills can be expensive. Improve register usage per thread by

  - Increasing the registers available per thread to 256

  - Running in lower SIMD width

- Passing large register file with online compilation.

```
ze_module_desc_t lz_mod_desc = {ZE_STRUCTURE_TYPE_MODULE_DESC,
                                nullptr,
                                ZE_MODULE_FORMAT_IL_SPIRV,
                                il_binary.size(),
                                il_binary.data(),
                                "-ze-opt-large-register-file",   // flags
                                nullptr};
```

- Enforcing SIMD16 with OpenCL kernel

  *__attribute__(intel_reqd_sub_group_size(16))*

  https://registry.khronos.org/OpenCL/extensions/intel/cl_intel_required_subgroup_size.html

Argonne
NATIONAL LABORATORY

# Where does performance stand currently

| Kernel | A100 | PVC (before optimization) | PVC (after optimizations) |
|---|---|---|---|
| IJacobian | 0.9 ms | 29 ms | 1.1 ms |
| Grad | 1.3 ms | 6.4 ms | 4.5 ms |
| Total time per timestep | 3.6 s | 70 s | 11 s |

*This is still a work in progress, and not a representative of performance between two hardware

Argonne
NATIONAL LABORATORY

# Key Takeaways

## Intel online compiler SYCL extension

- Allows for NVRTC-like runtime compilation
- Currently restricted to use of OpenCL C/LevelZero kernels (interop)
- In future, likely to extend support for SYCL kernel code

## Optimization Strategies

- Specialization constants - useful for optimizing code with parametric values. Runtime increases with number of spec. constants
- Use of appropriate workgroup sizes and barriers.
- For register-heavy kernels – large GRF + smaller SIMD width

## Future Work

- Systematic profiling of online compiled code for further insights on performance bottlenecks.
- Use of the updated SYCL kernel compiler with SYCL user functions.

Argonne
NATIONAL LABORATORY

# Acknowledgements

Team members:   Jed Brown (UC Boulder)

Natalie Beams (UT Knoxville)

Kenneth Jansen (UC Boulder)

**\*This work was done on a pre-production supercomputer with early versions of the Aurora software development kit.**

Argonne
NATIONAL LABORATORY