



LUMI, the Queen of the North

www.lumi-supercomputer.eu

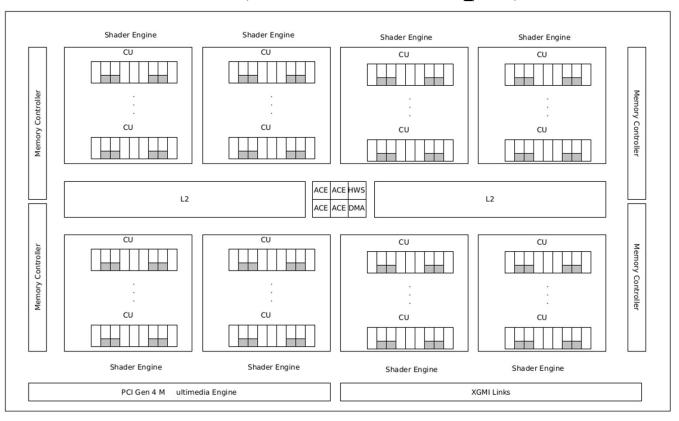
sharing and staging data

LUMI is a Tier-o GPU-accelerated Tier-o GPU partition: over LUMI-G: **supercomputer** that enables the 550 Pflop/s powered by **GPU** convergence of high-LUMI-D: AMD Instinct GPUs Partition CUMI-C: performance computing, Data x86 artificial intelligence, and high-Analytics Partition Partition performance data analytics. Interactive partition with 32 TB of memory and graphics Supplementary CPU LUMI-K: GPUs for data analytics and LUMI-F: partition Container High-speed Accelerated visualization Cloud ~200,000 AMD EPYC interconnect Storage Service CPU cores 7 PB Flash-based storage Possibility for combining layer with extreme I/O LUMI-Q: LUMI-P: different resources within a bandwidth of 2 TB/s and Lustre Emerging LUMI-O: single run. HPE Slingshot Storage IOPS capability. Cray tech Object technology. ClusterStor E1000. Storage Service 30 PB encrypted object storage (Ceph) for storing, 80 PB parallel file system

#lumisupercomputer #lumieurohpc

AMD GPUs (MI100 example)





LUMI will have the next generation of AMD instinct GPU

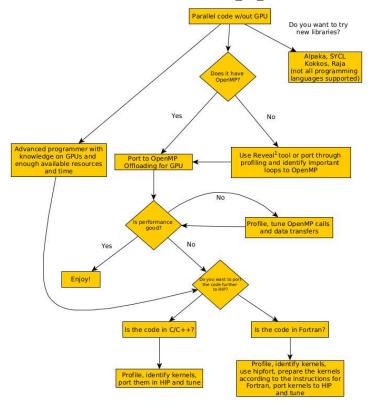


Differences between HIP and CUDA

- AMD GCN hardware wavefronts size is 64 (like warp for CUDA), some terminology is different
- Some CUDA library functions do not have AMD equivalents
- Shared memory and registers per thread can differ between AMD and NVIDIA hardware

LUMI

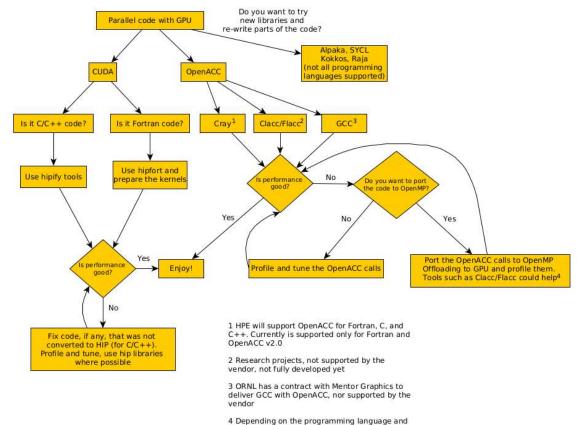
Porting Codes to LUMI (I)



1 Reveal will work good with Fortran codes and less with C, especially C++

LUMI

Porting Codes to LUMI (II)



if Clacc/Flacc can handle all the calls



BabelStream

- A memory bound benchmark from the university of Bristol
- Five kernels
 - o add (a[i]=b[i]+c[i])
 - o multiply (a[i]=b*c[i])
 - o copy (a[i]=b[i])
 - o triad (a[i]=b[i]+d*c[i])
 - o dot (sum = sum+d*c[i])



OpenMP Offloading

- Some basic OpenMP useful constructs:
 - o #pragma omp target enter/exit data map
 - o #pragma omp target teams distribute parallel for simd
 - o thread_limit(X) num_teams(Y)
- OpenMP 5.0, what is new: https://www.openmp.org/spec-html/5.0/openmpse71.html
- OpenMP 5.1, what is new: https://www.openmp.org/wp-content/uploads/OpenMP-API-Additional-Definitions-2-0.pdf



Improving performance on BabelStream for MI100

• Original call:

#pragma omp target teams distribute parallel for simd

• Optimized call

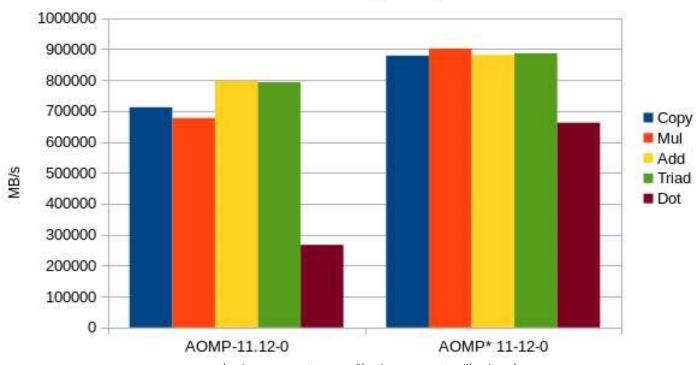
#pragma omp target teams distribute parallel for simd thread_limit(256) num_teams(240)

• For the dot kernel we used 720 teams



BabelStream, tune AOMP

BabelStream (MI100)



www.lumi-supercomputer.eu #lumisupercomputer #lumieurohpc



Introduction to HIP

- HIP: Heterogeneous Interface for Portability is developed by AMD to program on AMD GPUs
- It is a C++ runtime API and it supports both AMD and NVIDIA platforms
- HIP is similar to CUDA and there is no performance overhead on NVIDIA GPUs
- Many well-known libraries have been ported on HIP
- New projects or porting from CUDA, could be developed directly in HIP

https://github.com/ROCm-Developer-Tools/HIP



Hipify

- Hipify tools convert automatically CUDA codes
- It is possible that not all the code is converted, the remaining needs the implementation of the developer
- Hipify-perl: text-based search and replace
- Hipify-clang: source-to-source translator that uses clang compiler
- Porting guide: https://github.com/ROCm-Developer-Tools/HIP/blob/main/docs/markdown/hip-porting-guide.md



Hipify-perl

• It can scan directories and converts CUDA codes with replacement of the cuda to hip (sed –e 's/cuda/hip/g')

\$ hipify-perl --inplace filename

It modifies the filename input inplace, replacing input with hipified output, save backup in .prehip file.

\$ hipconvertinplace-perl.sh directory

It converts all the related files that are located inside the directory



Differences between CUDA and HIP API

CUDA

#include "cuda.h"

cudaMalloc(&d_x, N*sizeof(double));

cudaMemcpy(d_x,x,N*sizeof(double),
cudaMemcpyHostToDevice);

cudaDeviceSynchronize();

HIP

#include "hip/hip_runtime.h"
hipMalloc(&d_x, N*sizeof(double));

hipMemcpy(d_x,x,N*sizeof(double),
hipMemcpyHostToDevice);

hipDeviceSynchronize();



Launching kernel with CUDA and HIP

CUDA

HIP

```
hipLaunchKernelGGL(kernel_name,
gridsize,
blocksize,
shared_mem_size,
stream,
arg0, arg1, ...);
```



Libraries (not exhaustive)

NVIDIA	HIP	ROCm	Description
cuBLAS	hipBLAS	rocBLAS	Basic Linear Algebra Subroutines
cuRAND	hipRAND	rocRAND	Random Number Generator Library
cuFFT	hipFFT	rocFFT	Fast Fourier Transfer Library
cuSPARSE	hipSPARSE	rocSPARSE	Sparse BLAS + SPMV
NCCL		RCCL	Communications Primitives Library based on the MPI equivalents
CUB	hipCUB	rocPRIM	Low Level Optimized Parallel Primitives



Benchmark MatMul cuBLAS, hipBLAS

- Use the benchmark https://github.com/pc2/OMP-Offloading
- Matrix multiplication of 2048 x 2048, single precision
- All the CUDA calls were converted and it was linked with hipBlas
- CUDA (V100)

matMulAB (10): 1011.2 GFLOPS 12430.1 GFLOPS

• HIP (MI100)

matMulAB (10): 2327.6 GFLOPS 22216.7 GFLOPS

• MI100 achieves close to the theoretical peak for single precision



N-BODY SIMULATION

- N-Body Simulation (https://github.com/themathgeek13/N-Body-Simulations-CUDA) AllPairs N2
- 171 CUDA calls converted to HIP without issues, close to 1000 lines of code
- 32768 number of small particles, 2000 time steps

CUDA execution time on V100: 68.5 seconds

HIP execution time on MI100: 95.57 seconds, 39.5% worse performance

• Tune the number of threads per block to 256 instead of 1024, then:

HIP execution time on MI100: 54.32 seconds, 26.1% better performance than V100



Fortran

- First Scenario: Fortran + CUDA C/C++
 - OAssuming there is no CUDA code in the Fortran files.
 - oHipify CUDA
 - oCompile and link with hipco
- Second Scenario: CUDA Fortran
 - oThere is no HIP equivalent
 - oHIP functions are callable from C, using 'extern C'
 - oSee hipfort

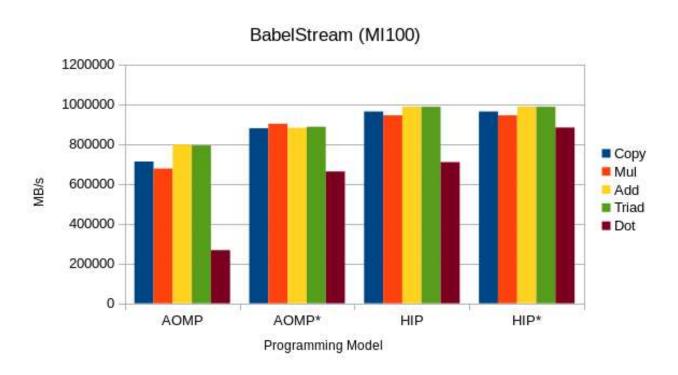


Hipfort

- The approach to port Fortran codes on AMD GPUs is different, the hipify tool does not support it.
- We need to use hipfort, a Fortran interface library for GPU kernel *
- Steps:
 - 1) We write the kernels in a new C++ file
 - 2) Wrap the kernel launch in a C function
 - 3) Use Fortran 2003 C binding to call the function
 - 4) Things **could** change in the future (see GPUFORT)
- Example of Fortran with HIP: https://github.com/cschpc/lumi/tree/main/hipfort
- Use OpenMP offload to GPUs
- * https://github.com/ROCmSoftwarePlatform/hipfort



BabelStream on MI100 (HIP vs AOMP)



www.lumi-supercomputer.eu #lumisupercomputer #lumieurohpc

LUMI

Megahip

- https://github.com/zjin-lcf/oneAPI-DirectProgramming
- 115 Applications/Examples with CUDA, SYCL, OpenMP offload and HIP
- Testing hipify tool, create a megahip script to convert all the CUDA examples to HIP
- ./megahip.sh

3287 CUDA calls were converted to HIP

115 applications totally 45692 lines of code, there are warnings for 4 of them, there are totally 24 warnings that something was wrong, check warnings.txt Application Success 96.5217% Conversion Success 99.2699%



OpenACC

- GNU will provide OpenACC (Mentor Graphics contract, now called Siemens EDA)
- HPE will use the provided GNU compiler for OpenACC support
- HPE is supporting for OpenACC v2.0 for Fortran. This is quite old OpenACC version. HPE announced support for OpenACC, newer versions for all the main programming languages (Fortran/C/C++)
- Clacc from ORNL: https://github.com/llvm-doe-org/llvm-project/tree/clacc/master OpenACC from LLVM only for C (Fortran and C++ in the future)
 - oTranslate OpenACC to OpenMP Offloading

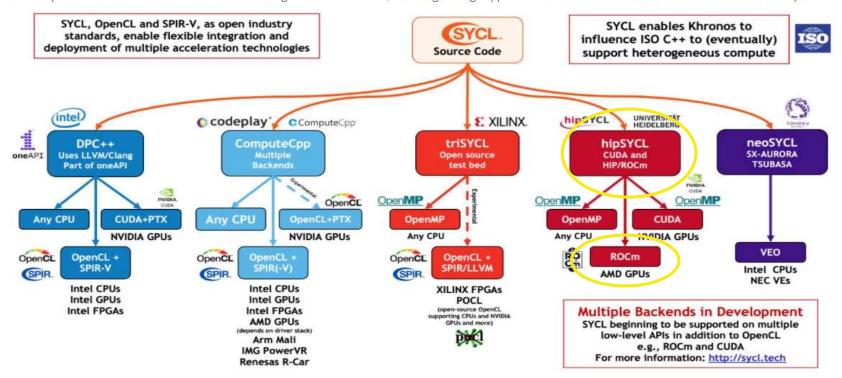


Clacc

- It supports C programming language, Fortran is on the way, C++ not started(??) yet \$ clang -fopenacc-print=omp -fopenacc-structured-ref-count-omp=no-hold -fopenacc-present-omp=no-present jacobi.c
- Original code (OpenACC): #pragma acc parallel loop reduction(max:lnorm) private(i,j) present(newarr, oldarr) collapse(2)

SYCL Implementations in Development

SYCL implementations are available from an increasing number of vendors, including adding support for diverse acceleration API back-ends in addition to OpenCL.



SAXPY SYCL

LUM

create queue

```
sycl::queue q(sycl::default_selector {});
const float A(aval);
sycl::buffer<float,1> d_X { h_X.data(), sycl::range<1>(h_X.size()) };
sycl::buffer<float,1> d_Y { h_Y.data(), sycl::range<1>(h_Y.size()) };
sycl::buffer<float,1> d_Z { h_Z.data(), sycl::range<1>(h_Z.size()) };
q.submit([&](sycl::handler& h) {
    auto X = d_X.template get_access<sycl::access::mode::read>(h);
    auto Y = d_Y.template get_access<sycl::access::mode::read>(h);
    auto Z = d_Z.template get_access<sycl::access::mode::read_write>(h);
h.parallel_for<class nstream>( sycl::range<1>{length}, [=] (sycl::id<1> i) {
        const int i = it[0];
        Z[i] = A * X[i] + Y[i];
    });
});
q.wait();
```

```
sycl::queue q(sycl::host_selector{});
sycl::queue q(sycl::cpu_selector{});
sycl::queue q(sycl::gpu_selector{});
sycl::queue q(sycl::accelerator_selector{});
```

Declare SYCL buffers to handle data on the device

SYCL accesors they generate a dataflow graph that the compiler and runtime can use to move data across devices

```
SYCL 2020
q.parallel_for( sycl::range<1>{length}, [=] (sycl::id<1> i) {
    d_Z[i] += A * d_X[i] + d_Y[i];
  });
```

Kokkos



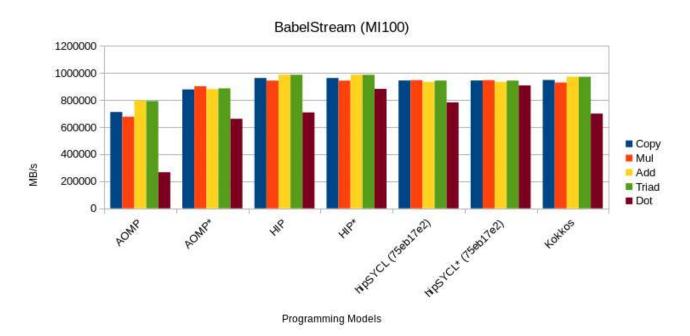
- A C++ programming framework
- Many concepts, such as view, execution pattern, execution policy, execution space, and computational body to support various devices
- It is used by a few HPC applications already with success
- There is some support for Fortran codes

```
Kokkos::View<T*> a(*d_a);
Kokkos::View<T*> b(*d_b);
...

Kokkos::parallel_reduce(array_size, KOKKOS_LAMBDA (const long index, T &tmp)
{
   tmp += a[index] * b[index];
}, sum);
```



Results of BabelStream on MI100



Kokkos is not optimized, compare with HIP results, not HIP*



Tuning

- Multiple wavefronts per compute unit (CU) is important to hide latency and instruction throughput
- Tune number of threads per block, number of teams for OpenMP offloading etc.
- Memory coalescing increases bandwidth
- Unrolling loops allow compiler to prefetch data
- Small kernels can cause latency overhead, adjust the workload
- Use of Local Data Share (LDS) memory



Conclusion/Future work

- A code written in C/C++ and MPI+OpenMP is a bit easier to be ported to OpenMP offloading compared to other approaches.
- The hipSYCL and Kokos could be a good option considering that the code is in C++.
- There can be challenges, depending on the code and what GPU functionalities are integrated to an application
- It will be required to tune the code for high occupancy
- Track historical performance among new compilers
- GCC for OpenACC and OpenMP Offloading for AMD GPUs (issues will be solved with GCC 12.x and LLVM 13.x)
- Tracking how profiling tools work on AMD GPUs (rocprof, TAU, HPCToolkit)
- We have trained more than 80 people on HIP porting: http://github.com/csc-training/hip



George Markomanolis

Lead HPC Scientist

georgios.markomanolis@csc.fi

Follow us

Twitter: @LUMIhpc

LinkedIn: <u>LUMI supercomputer</u>

YouTube: <u>LUMI supercomputer</u>

www.lumi-supercomputer.eu

contact@lumi-supercomputer.eu





The acquisition and operation of the EuroHPC supercomputer is funded jointly by the EuroHPC Joint Undertaking, through the European Union's Connecting Europe Facility and the Horizon 2020 research and innovation programme, as well as the of Participating States FI, BE, CH, CZ, DK, EE, IS, NO, PL, SE.







