

Performance Portability of the Chapel Language on Heterogeneous Architectures

Heterogeneity in Computing Workshop
IPDPS, 27 May 2024

Josh Milthorpe,
Xianghao Wang, and Ahmad Azizi

ORNL Advanced Computing Systems Research, milthorpejj@ornl.gov
Australian National University



ORNL is managed by UT-Battelle LLC for the US Department of Energy

This research used resources of the Experimental Computing Laboratory (ExCL) at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725.

Chapel Programming Language



- Parallel programming language supporting productive app development, including:
 - data exploration
 - multi-physics CFD
 - computational astrophysics
- Single-source compilation to multiple targets through LLVM
- First-class language features for task & data parallelism, synchronization, distributed memory
- Rapidly-improving GPU support
 - host-side code gen for memory management, kernel launch, synchronization
 - NVIDIA (LLVM PTX backend)
 - AMD (GCN backend)

Performance Portability

- Application code should run on many different hardware platforms ...
 - (without requiring rewriting for each new platform)
- ... and achieve acceptable performance on each platform
 - (without platform-specific optimizations)
- Pennycook, Sewall, and Lee's metric Φ : harmonic mean of efficiency on each platform
 - Architectural efficiency e.g. fraction of peak FLOP/s
 - Application efficiency e.g. inverse speedup versus fastest version
 - $\Phi = 0$ if code doesn't run on all platforms
- How well does Chapel support development of performance-portable application codes compared to more widely-used programming models like OpenMP and Kokkos?

S. J. Pennycook, J. D. Sewall, and V. W. Lee, [Implications of a metric for performance portability](#),
Future Generation Computer Systems, vol. 92, pp. 947–958, 2019.

Mini-apps

- We created new Chapel implementations of three mini-apps developed by the University of Bristol's High Performance Computing group
- These miniapps have been used extensively to compare parallel programming models and already have idiomatic implementations in OpenMP, Kokkos, CUDA, and HIP.
 - BabelStream: streaming memory access
 - miniBUDE: numerically intensive molecular dynamics
 - TeaLeaf: memory-intensive stencil PDE solver
- Not included in this study:
 - multi-device
 - distributed memory
 - programmer productivity

BabelStream

- An update of McCalpin's Stream memory bandwidth benchmark, comprising:

Kernel	Function	Load/Store	FLOP
Copy	$C = A$	2	0
Add	$C = A + B$	3	1
Mul	$B = \alpha * C$	2	1
Triad	$A = B + \alpha * C$	3	2 (1 FMA)
Nstream (PRK)	$A += B + \alpha * C$	4	3 (1 FMA)
Dot	$x = A . B$	2	2 (1 FMA)

- We measure BabelStream version 5.0 triad with 2^{28} 64-bit FP elements

<https://github.com/milthorpe/BabelStream>

Deakin, T., Price, J., Martineau, M., & McIntosh-Smith, S. (2018). [Evaluating attainable memory bandwidth of parallel programming models via BabelStream](#). International Journal of Computational Science and Engineering, 17(3), 247-262.

BabelStream Triad Implementations

- Chapel

```
proc triad() {  
  forall i in vectorDom do  
    A[i] = B[i] + scalar * C[i];  
  }  
}
```

 - CPU: loop is decomposed into chunks to be executed by worker threads
 - GPU: compiled to PTX (NVIDIA) or GCN (AMD) for each threads to compute a triad of elements; compiler generates host-side code for kernel launch and synchronization

```
const streamLocale = if useGPU  
  then here.gpus[deviceIndex]  
  else here;  
on streamLocale do {  
  const vectorDom = 0..#arraySize;  
  var A, B, C: [vectorDom] eltType = noinit;  
}
```

- CUDA

```
template <class T>  
void CUDASTream<T>::triad()  
{  
  triad_kernel<<<array_size/TBSIZE, TBSIZE>>>  
    (d_a, d_b, d_c);  
  check_error();  
  cudaDeviceSynchronize();  
  check_error();  
}
```

```
template <typename T>  
__global__ void triad_kernel(T * a, const T * b,  
const T * c)  
{  
  const T scalar = startScalar;  
  const int i = blockDim.x * blockIdx.x +  
threadIdx.x;  
  a[i] = b[i] + scalar * c[i];  
}
```

BabelStream Triad Implementations (2)

- Kokkos

```
template <class T>
void KokkosStream<T>::triad()
{
    Kokkos::View<T*> a(*d_a);
    Kokkos::View<T*> b(*d_b);
    Kokkos::View<T*> c(*d_c);

    const T scalar = startScalar;
    Kokkos::parallel_for(array_size, KOKKOS_LAMBDA (const long index)
    {
        a[index] = b[index] + scalar*c[index];
    });
    Kokkos::fence();
}
```

BabelStream Triad Implementations (3)

- OpenMP

```
template <class T>
void OMPStream<T>::triad()
{
    const T scalar = startScalar;

#ifdef OMP_TARGET_GPU
    int array_size = this->array_size;
    T *a = this->a;
    T *b = this->b;
    T *c = this->c;
    #pragma omp target teams distribute parallel for simd
#else
    #pragma omp parallel for
#endif
    for (int i = 0; i < array_size; i++)
    {
        a[i] = b[i] + scalar * c[i];
    }
    #if defined(OMP_TARGET_GPU) && defined(_CRAYC)
    // If using the Cray compiler, the kernels do not block, so this update forces
    // a small copy to ensure blocking so that timing is correct
    #pragma omp target update from(a[0:0])
    #endif
}
```

Experimental Platforms

	Processor	Sockets	Cores	Clock GHz	FP TFLOP/s	Mem BW GB/s	STREAM Balance*
CPU	Intel Skylake	2	8	3.70	1.89	256.0	59.2
	Intel Cascade Lake	2	24	4.00	6.14	287.3	171.1
	Intel Sapphire Rapids	2	52	3.80	12.65	614.4	164.7
	AMD Rome	2	64	3.00	6.14	409.6	120.0
	AMD Milan	2	32	3.68	3.77	409.6	73.6
	ARM ThunderX2	2	28	2.20	0.99	341.2	23.1
GPU	IBM POWER9	2	21	3.50	1.18	340.0	27.8
	NVIDIA P100	1	56	1.19	4.76	549.1	69.4
	NVIDIA V100	1	80	1.30	7.83	897.0	69.9
	NVIDIA A100	1	108	1.07	9.75	1935.0	40.3
	AMD MI60	1	64	1.20	7.37	1024.0	57.6
	AMD MI100	1	120	1.00	11.54	1229.0	75.1
	AMD MI250X	1	110	1.00	23.94	1600.0	119.7

* GFLOP s⁻¹ / Gword s⁻¹

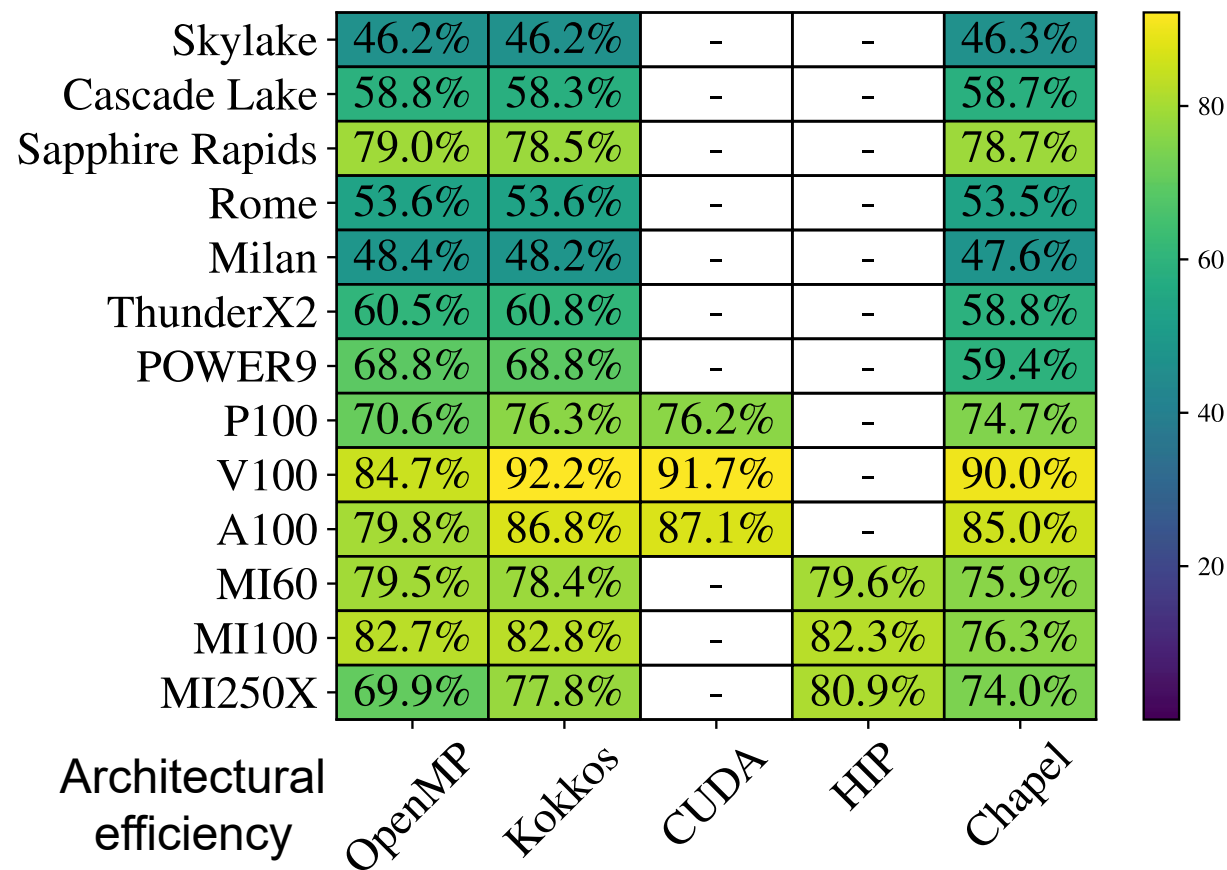
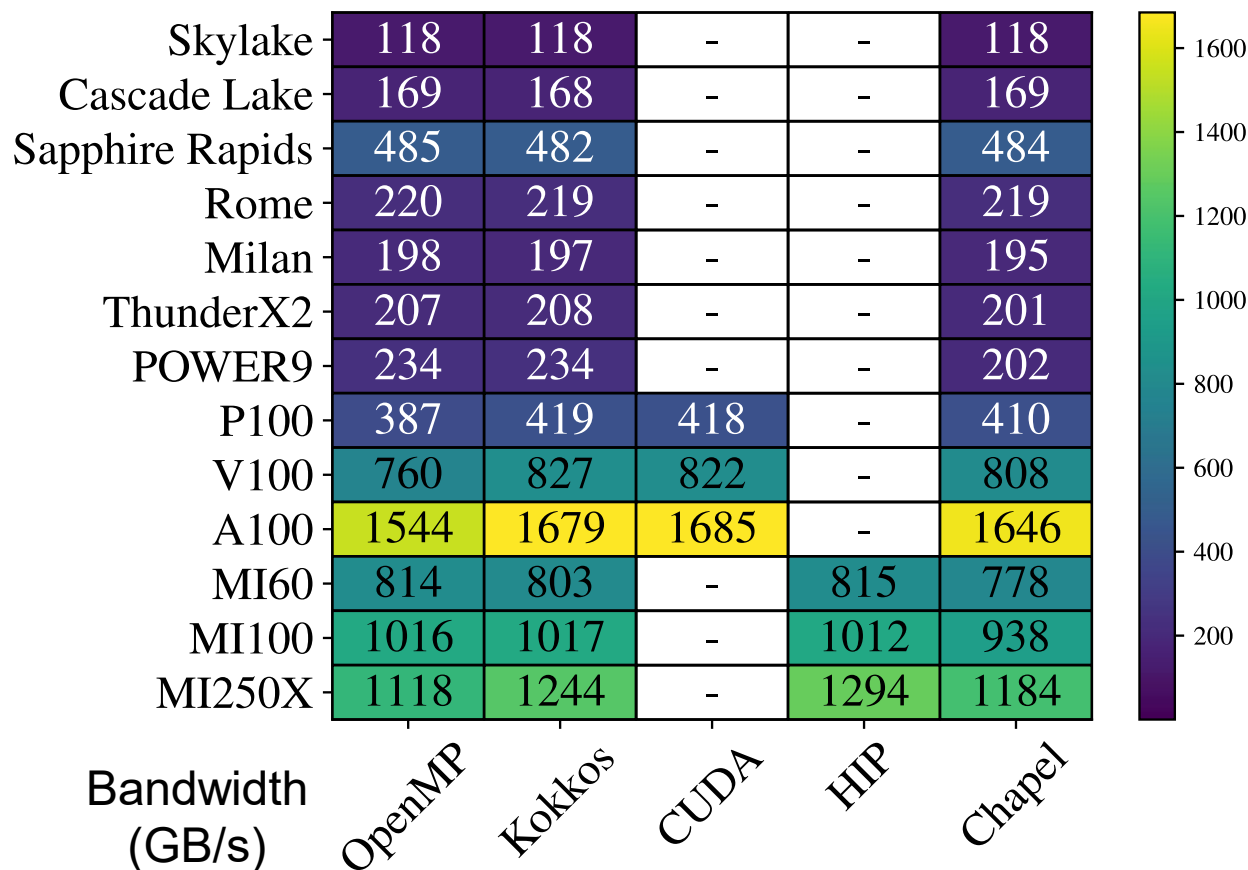
Experimental Configuration

	Processor	Operating System	GPU Driver Version	Compiler
CPU	Intel Skylake	Ubuntu 20.04.6		clang 17.0.6
	Intel Cascade Lake	Ubuntu 22.04.3		clang 17.0.1
	Intel Sapphire Rapids	Ubuntu 22.04.3		clang 17.0.1
	AMD Rome	Ubuntu 22.04.3		clang 17.0.6
	AMD Milan	Ubuntu 22.04.3		clang 17.0.6
	ARM ThunderX2	CentOS Stream 8		clang 17.0.2
	IBM POWER9	CentOS 8.3		gcc 10.2
GPU	NVIDIA P100	Ubuntu 20.04.6	525.147.05	nvcc 11.5
	NVIDIA V100	Ubuntu 22.04.3	550.54.15	nvcc 12.3
	NVIDIA A100	Ubuntu 22.04.3	555.42.02	nvcc 12.3
	AMD MI60	Ubuntu 22.04.3	6.3.6	hipcc 5.4.3
	AMD MI100	Ubuntu 22.04.3	5.15.0-15	hipcc 5.4.3
	AMD MI250X	SUSE LES 15.4	6.3.6	hipcc 5.4.3

Chapel 2.0, Kokkos 4.2.0

BabelStream Performance Portability

BabelStream v5 triad – 2²⁸ 64-bit elements



Platforms	OpenMP	Kokkos	CUDA	HIP	Chapel
All platforms	64.9%	65.8%	0	0	64.0%
Supported CPUs	57.5%	57.4%	0	0	56.1%
Supported GPUs	79.1%	82.9%	84.5%	80.9%	80.0%

miniBUDE

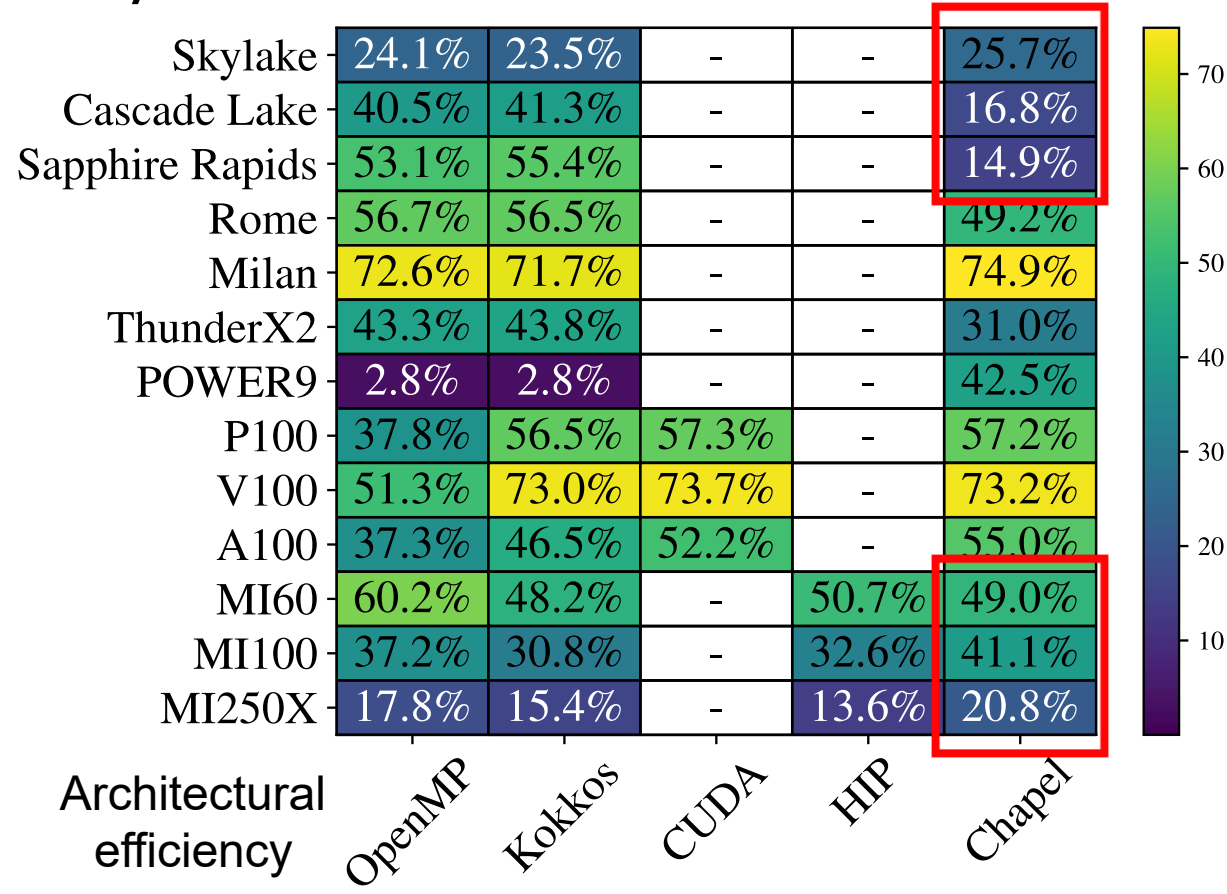
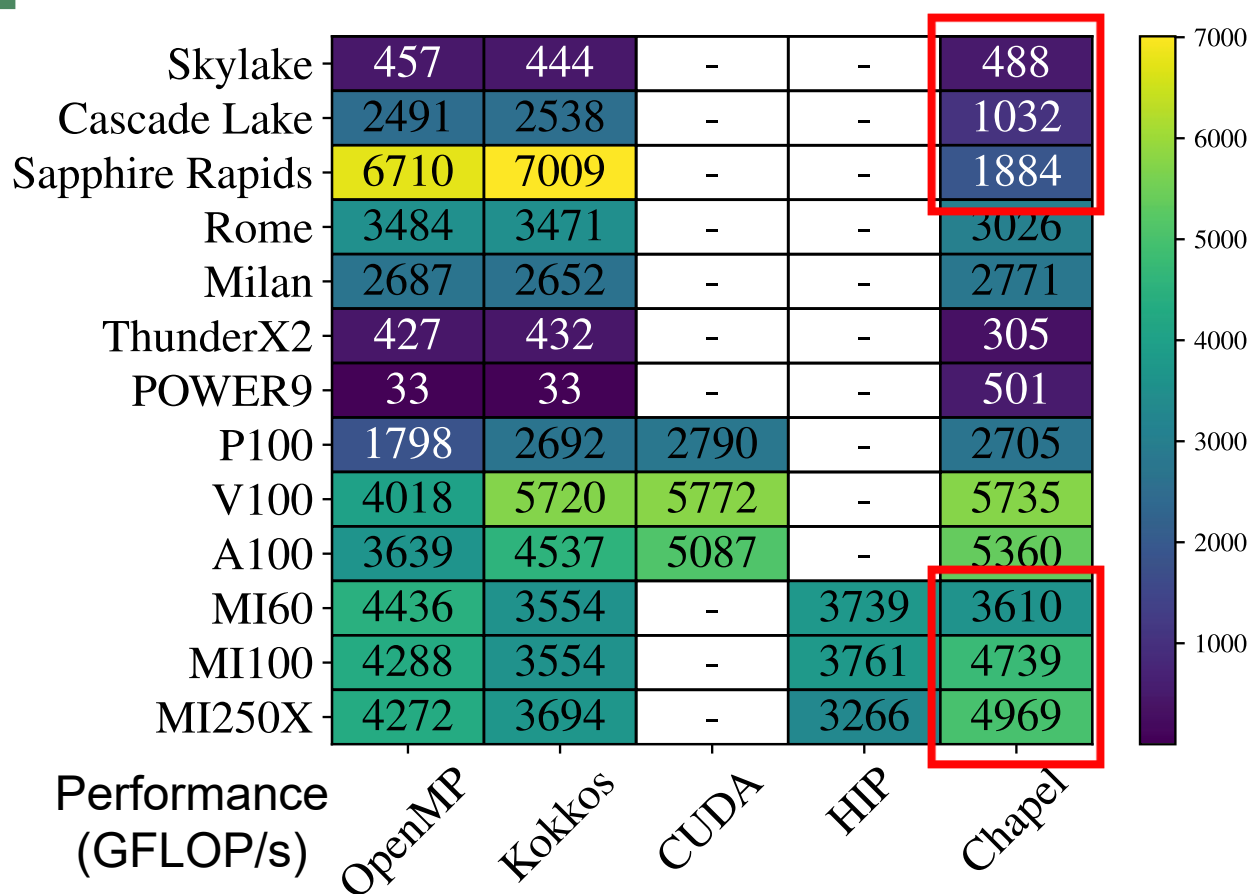
- Proxy app created from University of Bristol BUDE protein simulator
 - calculates energy of each ligand-protein pair in different poses (position + rotation)
 - highly arithmetically intensive: FP arithmetic & trigonometric
- Kernel is triply-nested loop over proteins, ligands, poses
 - Chapel follows CUDA decomposition: 1D kernel assigning multiple poses to thread

```
foreach group in 0..
```

```
const protein = context.protein;
const ligand = context.ligand;
const forcefield = context.forcefield;
const poses: [0:int(32)..<6:int(32), 0..
```

miniBUDE Performance Portability

miniBUDE v2 – small 'bm1' input



Platforms	OpenMP	Kokkos	CUDA	HIP	Chapel
All platforms*	43.0%	44.8%	0	0	33.8%
Supported CPUs*	43.0%	43.1%	0	0	25.9%
Supported GPUs	43.0%	47.1%	60.2%	39.7%	53.1%

* Except POWER9

TeaLeaf

- Collection of iterative sparse linear solvers, simulating heat conduction over time using five-point stencils over 2D grid
- Low arithmetic intensity = better suited to low STREAM balance
- 2D index domains: expose parallelism over both loops

```
#pragma omp target teams distribute parallel for simd
collapse(2)
for (int jj = halo_depth; jj < y - halo_depth; ++jj) {
    for (int kk = halo_depth; kk < x - halo_depth; ++kk) {
        const int index = kk + jj * x;
        p[index] = beta * p[index] + r[index];
    }
}
```

```
Kokkos::parallel_for(
    x * y, KOKKOS_LAMBDA(const int &index) {
        const int kk = index % x;
        const int jj = index / x;

        if (kk >= halo_depth
            && kk < x - halo_depth
            && jj >= halo_depth && jj < y - halo_depth) {
            p(index) = beta * p(index) + r(index);
        }
    });
```

```
[(i,j) in Domain.expand(-halo_depth)] p[i,j] = beta * p[i,j] + r[i,j];
```

<https://github.com/milthorpe/TeaLeaf>

S. McIntosh-Smith, et al., [Tealeaf: A mini-application to enable design-space explorations for iterative sparse linear solvers](#).
IEEE International Conference on Cluster Computing (CLUSTER), 2017.

TeaLeaf - Reductions

- Many sum reductions to compute global deltas or error metrics
 - In Chapel 2.0, these must be computed in global memory

```
Kokkos::parallel_reduce(
  x * y,
  KOKKOS_LAMBDA(const int &index, double &rrn_temp) {
    const int kk = index % x;
    const int jj = index / x;
    if (kk >= halo_depth
        && kk < x - halo_depth
        && jj >= halo_depth
        && jj < y - halo_depth) {
      u(index) += alpha * p(index);
      r(index) -= alpha * w(index);
      rrn_temp += r(index) * r(index);
    }
  },
  *rrn);
```

```
var temp: [reduced_local_domain] real = noinit;
...
forall oneDIdx in reduced_OneD {
  const ij = reduced_local_domain.orderToIndex(oneDIdx);
  u[ij] += alpha * p[ij];
  r[ij] -= alpha * w[ij];
  temp[ij] = r[ij] ** 2;
}
rrn = gpuSumReduce(temp);
```

Chapel 2.0

```
var rrn: real;
forall ij in reduced_local_domain
  with (+ reduce rrn) {
    u[ij] += alpha * p[ij];
    r[ij] -= alpha * w[ij];
    rrn += r[ij] ** 2;
  }
```

Chapel 2.x

TeaLeaf – Chapel Multi-Dimensional Indexing

- Using 2D indices improved readability of Chapel code and performed well on CPU platforms
- However, using 2D domains reduced GPU performance due to under-utilization of available GPU cores in Chapel 2.0
 - first dimension is assigned to GPU threads
 - remaining dimensions implemented as loops inside GPU kernel
- We replaced multi-dimensional loops with 1D loop over linearized space to allow full utilization of GPU cores

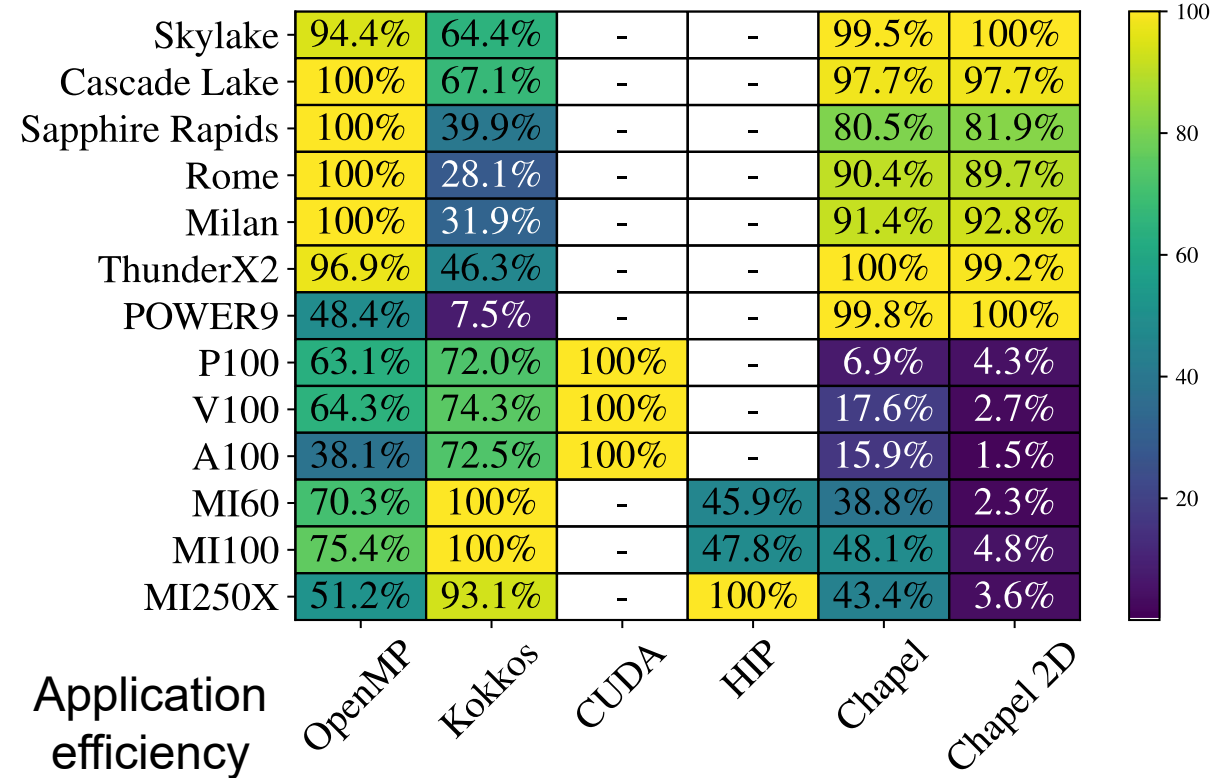
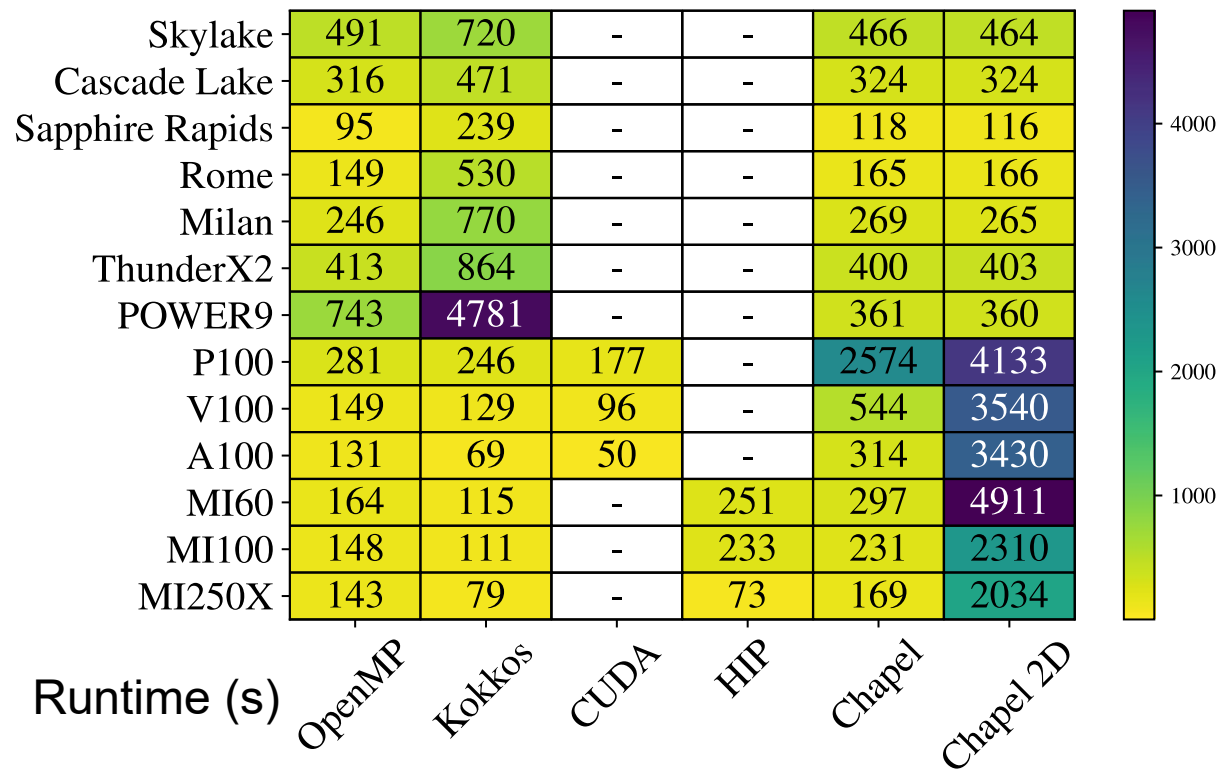
```
const Domain = {0..forall ij in Domain {  
  u[ij] = energy[ij] * density[ij];  
}
```



```
const Domain = {0..const OneD = {0..foreach oneDIdx in OneD {  
  const ij = local_domain.orderToIndex(oneDIdx);  
  u[ij] = energy[ij] * density[ij];  
}
```

TeaLeaf Performance Portability

tea_bm_5.in – 4000×4000 CG solve, 10 iters



Platforms	OpenMP	Kokkos	CUDA	HIP	Chapel	Chapel 2D
All platforms	69.8%	37.3%	0	0	31.5%	5.7%
Supported CPUs	85.8%	25.3%	0	0	93.7%	94.0%
Supported GPUs	57.3%	83.5%	100.0%	56.9%	17.8%	2.7%

Conclusions

- Pennycook, Sewall and Lee's metric Φ remains a useful lens for evaluating portable programming models and identifying areas of strength and weakness
- Performance portability of OpenMP and Kokkos continues to improve
- Chapel is a new option for performance-portable parallel programming
 - concise code
 - (mostly) good performance across a wide range of platforms
 - easier path to multi-device, multi-node distribution
- Some issues remain with Chapel GPU code generation
 - fixing these will avoid performance pitfalls for users

Acknowledgments

- This research used resources of the Experimental Computing Laboratory (ExCL) and the Oak Ridge Leadership Computing Facility at Oak Ridge National Laboratory, which are supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725, and resources of the National Computational Infrastructure (NCI), which is supported by the Australian Government.
- This research was funded, in part, by the Brisbane Advanced Programming Systems Project.
- Thanks to Engin Kayraklioglu, Brad Chamberlain, Michael Ferguson, and Jade Abraham from the HPE Chapel team for helpful discussions during our development of the mini-applications in Chapel.