

Exploring a multi-resolution GPU programming model for Chapel

Akihiro Hayashi (Georgia Tech)

Sri Raj Paul (Georgia Tech)

Vivek Sarkar (Georgia Tech)

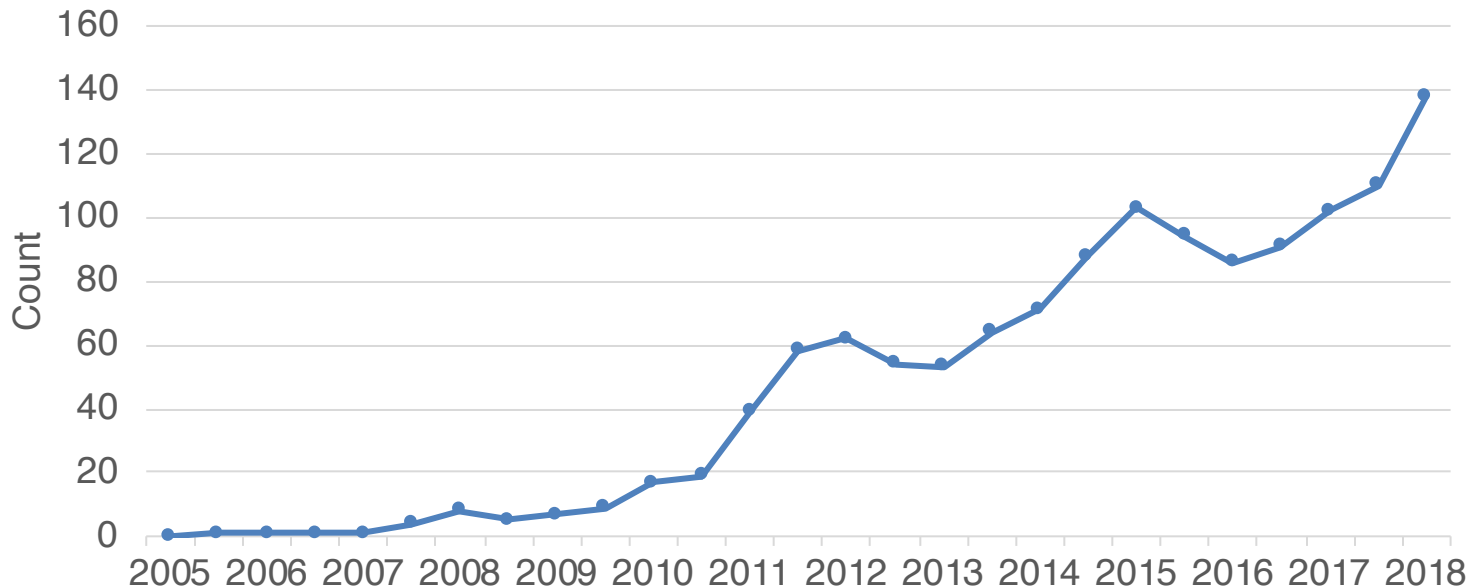


A multi-resolution GPU programming model for Chapel

MOTIVATION

GPUs are a common source of performance improvement in HPC

Accelerator/Co-Processor in Top500



**Aurora (ANL)
El Capitan(LLNL)
Frontier (ORNL)
all plan
to include GPUs!**

2021

Source: <https://www.top500.org/statistics/list/>



GPU Programming in Chapel

□ Chapel's multi-resolution concept

- High-level
- Start with writing “forall” loops (on CPU, proof-of-concept)

```
forall i in 1..n {  
    ...  
}
```

- Low-level
- Apply automatic GPU code generators [1][2] when/where possible
 - Consider using the GPUlterator module [3] (Mason-registry)
 - Consider writing GPU kernels using CUDA/HIP/OpenCL or other accelerator language, and invoke them from Chapel

[1] Albert Sidelnik et al. Performance Portability with the Chapel Language (IPDPS '12).

[2] Michael L. Chu et al. GPGPU support in Chapel with the Radeon Open Compute Platform (CHIUIW'17).

[3] Akihiro Hayashi et al. GPUlterator: Bridging the gap between Chapel and GPU Platforms (CHIUIW'19).



Example: STREAM (Original)

```
1 var A: [1..n] real(32);
2 var B: [1..n] real(32);
3 var C: [1..n] real(32);
4 // STREAM
5 forall i in 1..n {
6     A(i) = B(i) + alpha * C(i);
7 }
```



Example: STREAM (C interoperability)

- ❑ Invoking CUDA/HIP/OpenCL code using the C interoperability feature

```
1 extern proc GPUST(A: [] real(32),      1 // separate C file
2         B: [] real(32),                2 void GPUST(float *A,
3         C: [] real(32),                3         float *B,
4         a1: real(32),                  4         float *C,
5         lo: int, hi: int);             5         float a1,
6 var A: [1..n] real(32);                6         int start,
7 var B: [1..n] real(32);                7         int end) {
8 var C: [1..n] real(32);                8 // CUDA/HIP/
9 // Invoking CUDA/OpenCL program       9         OpenCL Code
10 GPUST(A, B, C, alpha, 1, n);          10 }
```



Example: STREAM (GPUIterator)

- ❑ Connecting the GPU version with the forall loop using the GPUIterator module

```
1 // GPU Iterator (in-between)
2 var G = lambda (lo: int, hi: int,
3               nElems: int) {
4     GPUST(A, B, C, alpha, 1, n);
5 };
6 var CPUPercent = 50;
7 forall i in GPU(1..n, G, CPUPercent) {
8     A(i) = B(i) + alpha * C(i);
9 }
```

```
1 // separate C file
2 void GPUST(float *A,
3           float *B,
4           float *C,
5           float a1,
6           int start,
7           int end) {
8     // CUDA/HIP/
9     OpenCL Code
10 }
```

Note: the GPUIterator is designed to facilitate

1) hybrid execution (CPUs+GPUs), and 2) distributed execution

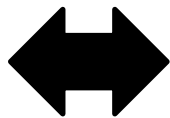


Example:

No appropriate GPU abstraction

Highest-level Chapel-GPU Programming

```
1 forall i in 1..n {  
2   A(i) = B(i) + alpha * C(i);  
3 }
```



A huge gap!

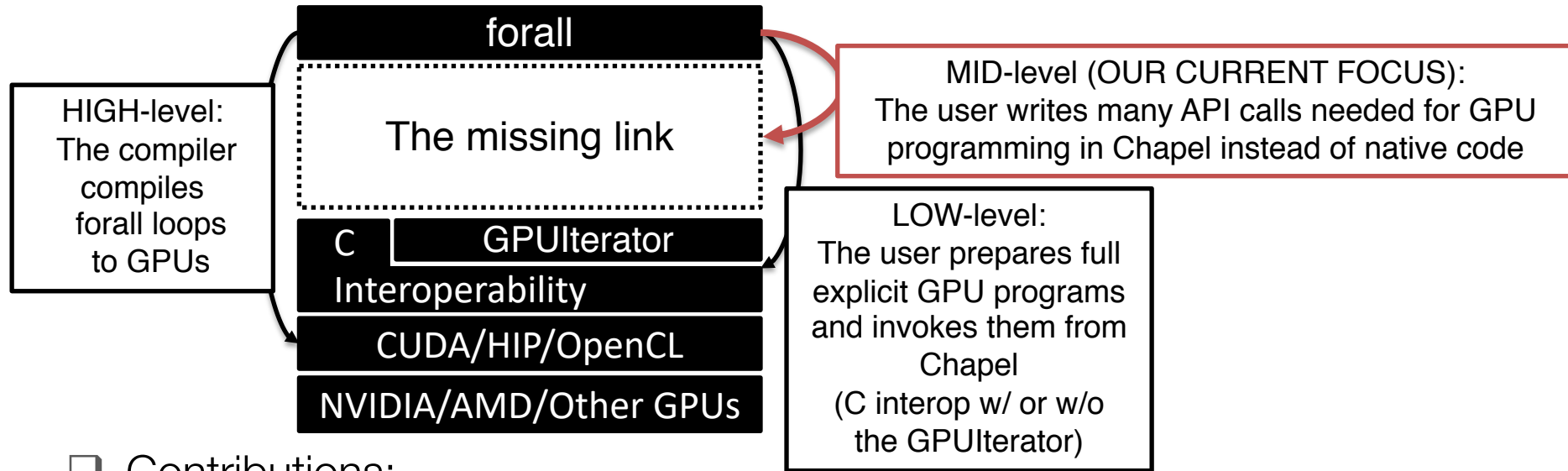
Lowest-level Chapel-GPU Programming
(C Interoperability only or GPUIterator)

```
1 // separate C file  
2 __global__ void stream(float *dA, float *dB, float *dC,  
3                       float alpha, int N) {  
4     int id = blockIdx.x * blockDim.x + threadIdx.x;  
5     if (id < N) {  
6         dA[id] = dB[id] + alpha * dC[id];  
7     }  
8 }  
9 void GPUST(float *A, float *B, float *C, float alpha,  
10            int start, int end, int GPUN) {  
11     float *dA, *dB, *dC;  
12     CudaSafeCall(cudaMalloc(&dA, sizeof(float) * GPUN));  
13     CudaSafeCall(cudaMalloc(&dB, sizeof(float) * GPUN));  
14     CudaSafeCall(cudaMalloc(&dC, sizeof(float) * GPUN));  
15     CudaSafeCall(cudaMemcpy(dB, B + start, sizeof(float) *  
16                          GPUN, cudaMemcpyHostToDevice));  
17     CudaSafeCall(cudaMemcpy(dC, C + start, sizeof(float) *  
18                          GPUN, cudaMemcpyHostToDevice));  
19  
20     stream<<<ceil(((float)GPUN)/1024), 1024>>>  
21                          (dA, dB, dC, alpha, GPUN);  
22     CudaSafeCall(cudaDeviceSynchronize());  
23     CudaSafeCall(cudaMemcpy(A + start, dA, sizeof(float) *  
24                          GPUN, cudaMemcpyDeviceToHost));  
25     CudaSafeCall(cudaFree(dA));  
26     CudaSafeCall(cudaFree(dB));  
27     CudaSafeCall(cudaFree(dC));  
28 }
```

Research Question:

What is an appropriate and portable programming interface that bridges the "forall" and GPU versions?

Big Picture: A Multi-level Chapel GPU Programming Model



□ Contributions:

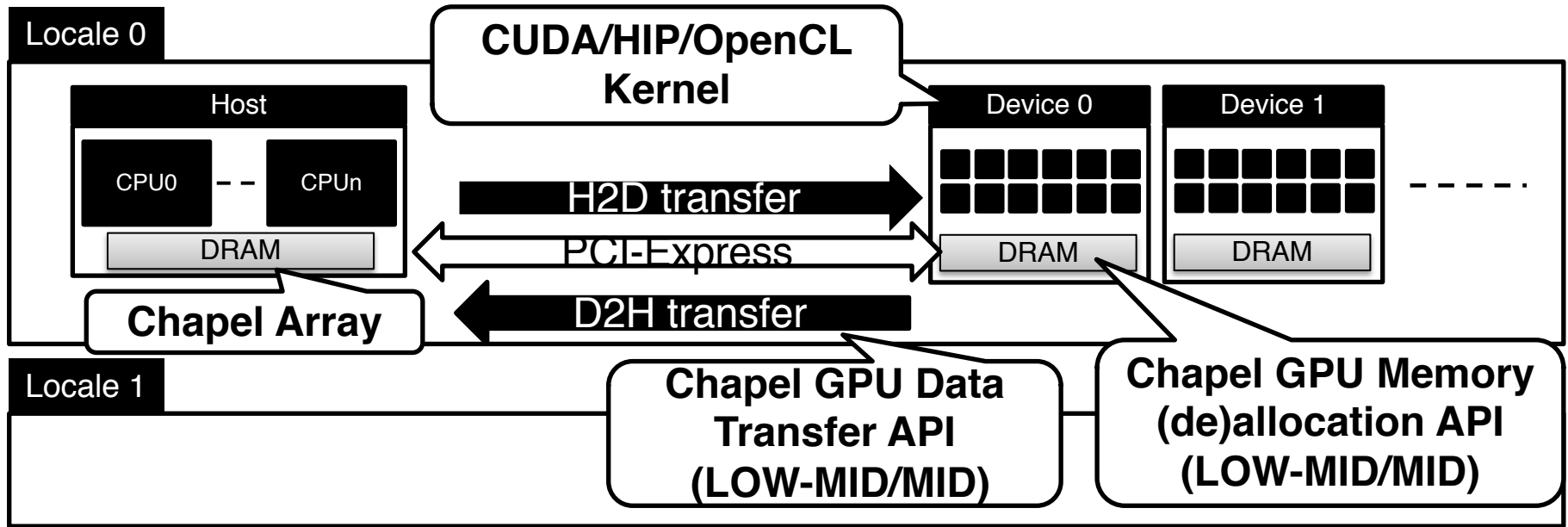
- Design and implementation of “MID”/” LOW-MID” levels Chapel GPU API
- Performance evaluations on different distributed and single-node platforms (Summit, Cori, and a single-node AMD machine)



A multi-resolution GPU programming model for Chapel

DESIGN

An overview of the LOW-MID/MID GPU API for Chapel



Note: Multi-locales plus multi-GPUs execution can be easily done with the GPUlterator module



Chapel GPU API Design: Summary

- ❑ Use case:
 - The user would like to 1) write GPU kernels, or 2) utilize highly-tuned GPU libraries, and would like to stick with Chapel for the other parts (allocation, data transfers)
- ❑ Provides two levels of GPU API
 - LOW-MID: Provides wrapper functions for raw GPU APIs
Example: `var ga: c_void_ptr = GPUAPI.Malloc(sizeInBytes);`
 - MID: Provides more user-friendly APIs
Example: `var ga = new GPUArray(A);`
- ❑ Important design decisions
 - The user is still supposed to write kernels in CUDA/HIP/OpenCL
 - The APIs significantly facilitates the orchestration of:
 - ✓ Device memory (de)allocation, and host-to-device/device-to-host data transfers,
 - The use of the APIs does not involve any modifications to the Chapel compiler



Chapel GPU API Design:

LOW-MID GPU API

□ Summary

- Provides the same functionality as CUDA/HIP/OpenCL
- The user is still supposed to write CUDA/HIP/OpenCL kernels
- The user is supposed to handle both C types and Chapel types

□ Key APIs

- Device Memory Allocation
 - ✓ `Malloc(ref devPtr: c_void_ptr, size: size_t);`
- Host-to-device, and device-to-host data transfers
 - ✓ `Memcpy(dst: c_void_ptr, src: c_void_ptr, count: size_t, direction: int);`
- Ensuring the completion of GPU computations
 - ✓ `DeviceSynchronize(void);`
- Device Memory deallocation
 - ✓ `Free(c_void_ptr);`



Chapel GPU API Design:

MID GPU API

□ Summary

- More natural to Chapel programmers
- The user is still supposed to write CUDA/HIP/OpenCL kernels

□ Key APIs

- Device Memory Allocation
 - ✓ `var dA = new GPUArray(A);`
- Host-to-device, and device-to-host data transfers
 - ✓ `ToDevice(dA:GPUArray, ...); FromDevice(dA: GPUArray, ...);`
 - ✓ `dA.ToDevice(); dA.fromDevice();`
- Device Memory deallocation
 - ✓ `Free(dA:GPUArray, ...);`
 - ✓ `dA.Free();`



Chapel GPU API Design: LOW-MID/MID GPU API Example

LOW-MID Level

```
1 use GPUAPI;
2 var A: [1..n] real(32);
3 var B: [1..n] real(32);
4 var C: [1..n] real(32);
5 var dA, dB, dC: c_void_ptr;
6 var size: size_t =
7     (A.size:size_t * c_sizeof(A.eltType));
8 Malloc(dA, size);
9 Malloc(dB, size);
10 Malloc(dC, size);
11 Memcpy(dB, c_ptrTo(B), size, TODVICE);
12 Memcpy(dC, c_ptrTo(C), size, TODVICE);
13 LaunchST(dA, dB, dC, alpha, N: size_t);
14 DeviceSynchronize();
15 Memcpy(c_ptrTo(A), dA, size, FROMDEVICE);
16 Free(dA); Free(dB); Free(dC);
```

MID-level

```
1 use GPUAPI;
2 var A: [1..n] real(32);
3 var B: [1..n] real(32);
4 var C: [1..n] real(32);
5 var dA = new GPUArray(A);
6 var dB = new GPUArray(B);
7 var dC = new GPUArray(C);
8 toDevice(dB, dC);
9 LaunchST(dA.dPtr(), dB.dPtr(),
10         dC.dPtr(), alpha,
11         dN: size_t);
12 DeviceSynchronize();
13 FromDevice(dA);
14 Free(dA, dB, dC);
```



Example: Single-node execution of STREAM (MID-level w/ GPUIterator)

```
1
2 var A: [1..n] real(32);
3 var B: [1..n] real(32);
4 var C: [1..n] real(32);
5 var GPUCallback = lambda (lo: int, hi: int, nElems: int) {
6     var dA = new GPUArray(A);
7     var dB = new GPUArray(B);
8     var dC = new GPUArray(C);
9     toDevice(dB, dC);
10    LaunchST(dA.dPtr(), dB.dPtr(),
11            dC.dPtr(), alpha,
12            dN: size_t);
13    DeviceSynchronize();
14    FromDevice(dA);
15    Free(dA, dB, dC);
16 };
17 forall i in GPU(1..n, GPUCallback,
18                CPUPercent) {
19     A(i) = B(i) + alpha * C(i);
20 }
```

```
1 // separate C file (CUDA w/ device lambda)
2 void LaunchST(float *dA, float *dB,
3              float *dC, float alpha, int N) {
4     GPU_FUNCTOR(N, 1024, NULL,
5                 [=] __device__ (int i) {
6         dA[i] = dB[i] + alpha * dC[i];
7     });
8 }
```

The user has the option of writing a device function(s), a device lambda(s), or a library call(s)

Example: Distributed execution of STREAM (MID-level w/ GPUIterator)

```
1 var D: domain(1) dmapped Block(boundingBox={1..n}) = {1..n};
2 var A: [D] real(32);
3 var B: [D] real(32);
4 var C: [D] real(32);
5 var GPUCallback = lambda (lo: int, hi: int, nElems: int) {
6     var dA = new GPUArray(A.localSlice(lo..hi));
7     var dB = new GPUArray(B.localSlice(lo..hi));
8     var dC = new GPUArray(C.localSlice(lo..hi));
9     toDevice(dB, dC);
10    LaunchST(dA.dPtr(), dB.dPtr(),
11            dC.dPtr(), alpha,
12            dN: size_t);
13    DeviceSynchronize();
14    FromDevice(dA);
15    Free(dA, dB, dC);
16 };
17 forall i in GPU(D, GPUCallback,
18                CPUPercent) {
19     A(i) = B(i) + alpha * C(i);
20 }
```

The user has the option of writing a device function(s), a device lambda(s), or a library call(s)

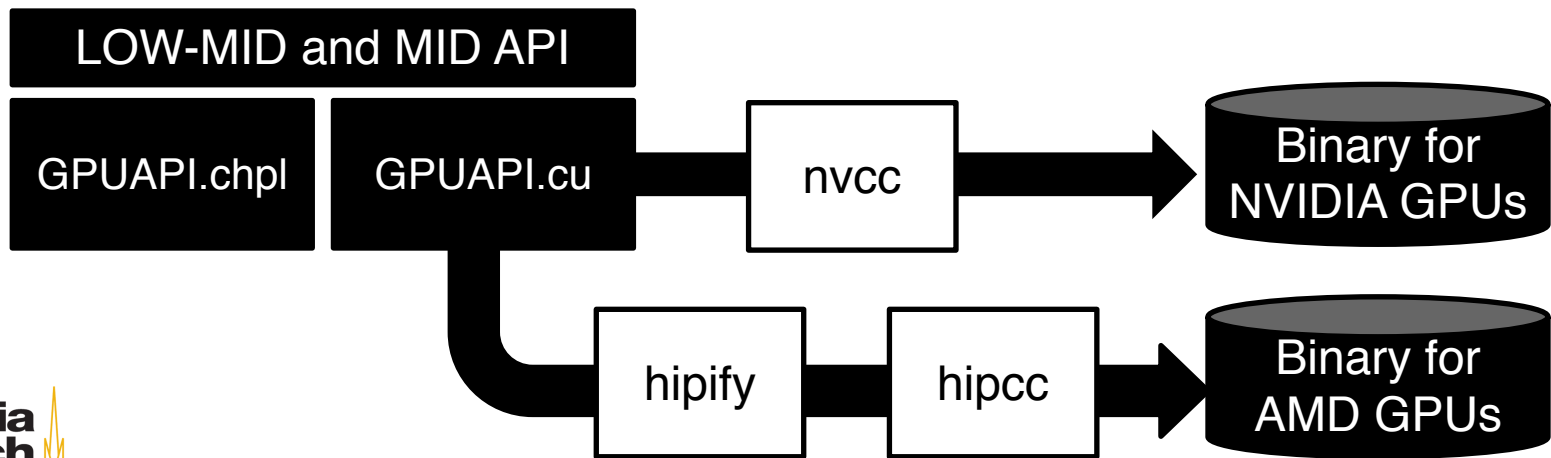
```
1 // separate C file (CUDA w/ device lambda)
2 void LaunchST(float *dA, float *dB,
3              float *dC, float alpha, int N) {
4     GPU_FUNCTOR(N, 1024, NULL,
5                 [=] __device__ (int i) {
6         dA[i] = dB[i] + alpha * dC[i];
7     });
8 }
```

A multi-resolution GPU programming model for Chapel

IMPLEMENTATION

Implementation of GPU API

- ❑ Provides an external module (GPUAPI)
 - Can be used either stand-alone or with the GPUIterator module
 - <https://github.com/ahayashi/chapel-gpu>
 - ✓ The “feature/explicit” branch
- ❑ Currently supports NVIDIA and ROCM-ready AMD GPUs



A multi-resolution GPU programming model for Chapel

PERFORMANCE EVALUATIONS

Performance Evaluations

□ Platforms

- Cori GPU@NERSC: Intel Xeon (Skylake) + NVIDIA V100 GPU
- Summit@ORNL: IBM POWER9 + NVIDIA Tesla V100 GPU
- A single-node AMD machine: Ryzen9 3900 + Radeon RX570

□ Chapel Compilers & Options

- Chapel Compiler 1.20.0 with the --fast option

□ GPU Compilers

- CUDA: NVCC 10.2 (Cori), 10.1 (Summit) with the -O3 option
- AMD: ROCM 2.9.6, HIPCC 2.8 with the -O3 option



Performance Evaluations (Cont'd)

- ❑ Tasking & Multi-locale execution
 - CHPL_TASK=qthreads
 - CHPL_COMM=gasnet
 - CHPL_COMM_SUBSTRATE=ibv
- ❑ GPUlterator (For distributed GPU execution)
 - <https://github.com/ahayashi/chapel-gpu/tree/feature/explicit>
- ❑ Applications
 - Stream
 - BlackScholes
 - Matrix Multiplicaiton
 - Logistic Regression
 - Source code can be found at:
<https://github.com/ahayashi/chapel-gpu/tree/feature/explicit/apps>



How many lines are added/modified?

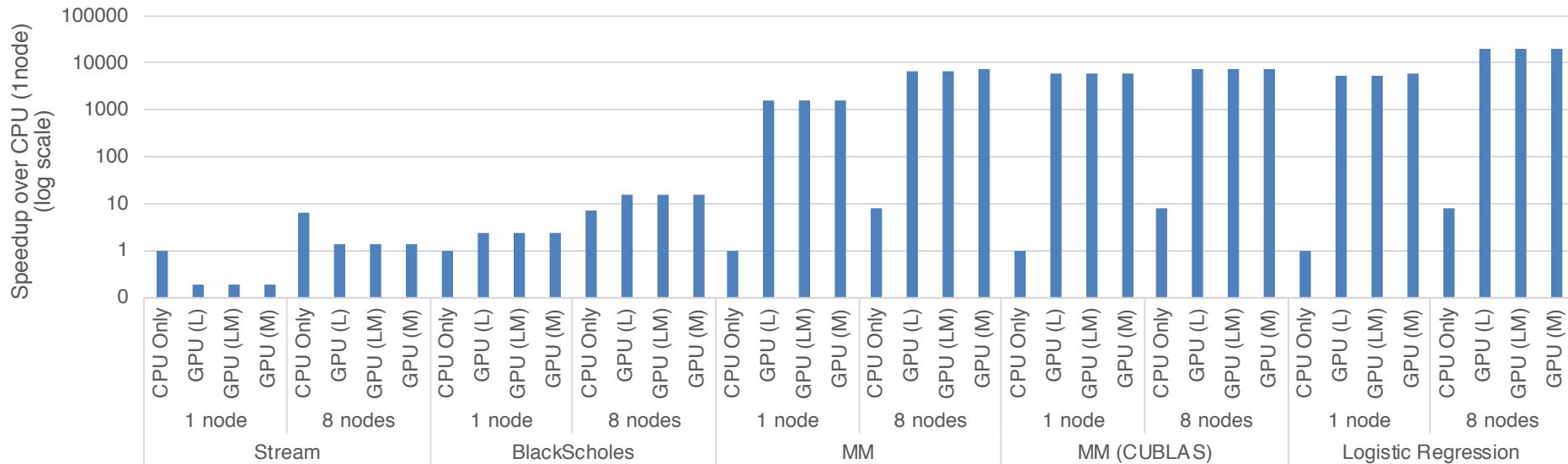
	LOC Baseline (Chapel LOW + CUDA)	LOC LOW-MID (Chapel LOW- MID + CUDA)	LOC MID (Chapel MID + CUDA)
Stream	4 + 24 = 28	16 + 9 = 25	11 + 9 = 20
BlackScholes	4 + 99 = 103	16 + 83 = 99	11 + 83 = 94
Logistic Regression	2 + 36 = 38	16 + 18 = 34	10 + 18 = 28
Matrix Multiplication	3 + 30 = 33	14 + 15 = 29	10 + 16 = 26

□ The use of GPU API decreases LOC



How fast are GPUs?

(Multi-nodes, 1 GPU/node, Cori)

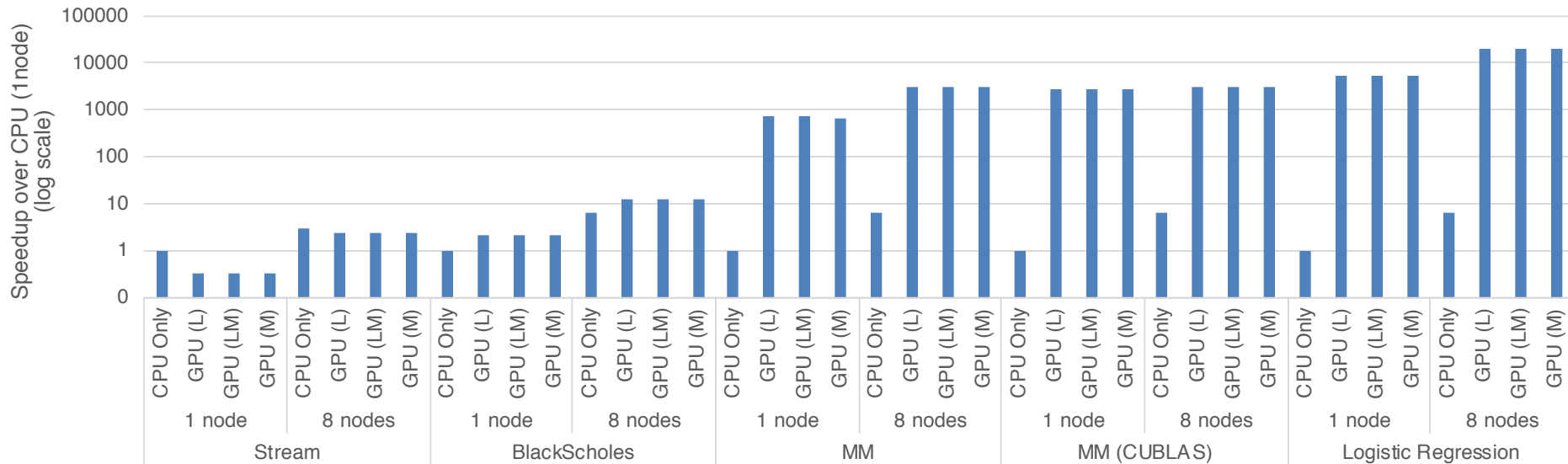


- Each app shows good strong scalability
- No significant performance difference between L (LOW), LM (LOW-MID), and M (MID)



How fast are GPUs?

(Multi-nodes, 1 GPU/node, Summit)



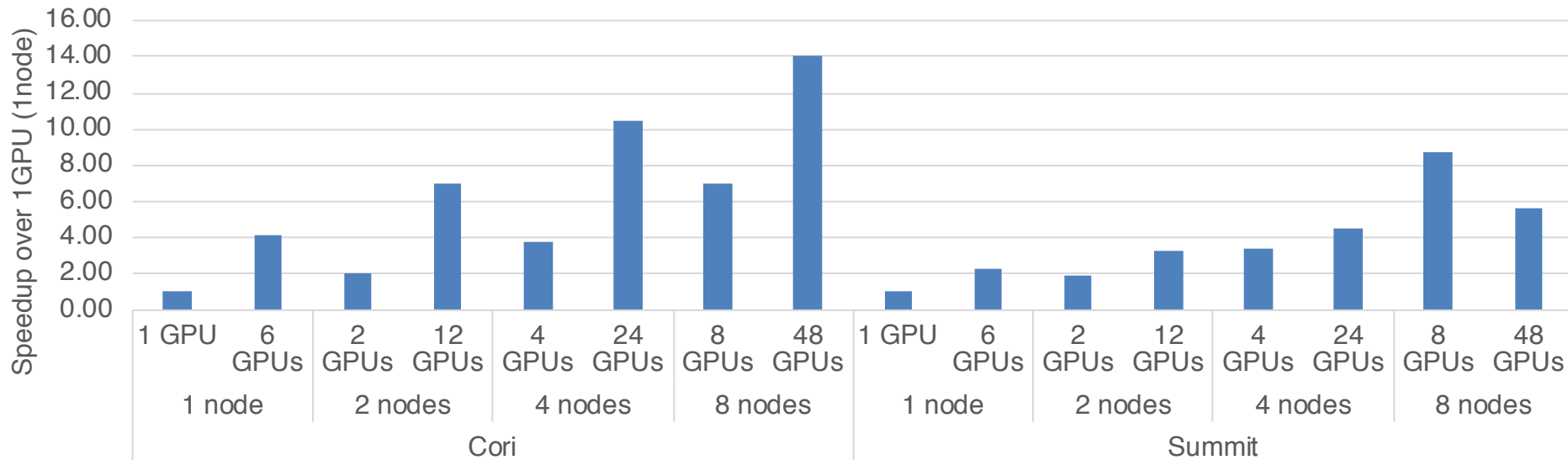
- ❑ Each app shows good strong scalability
- ❑ No significant performance difference between L (LOW), LM (LOW-MID), and M (MID)



How fast are GPUs?

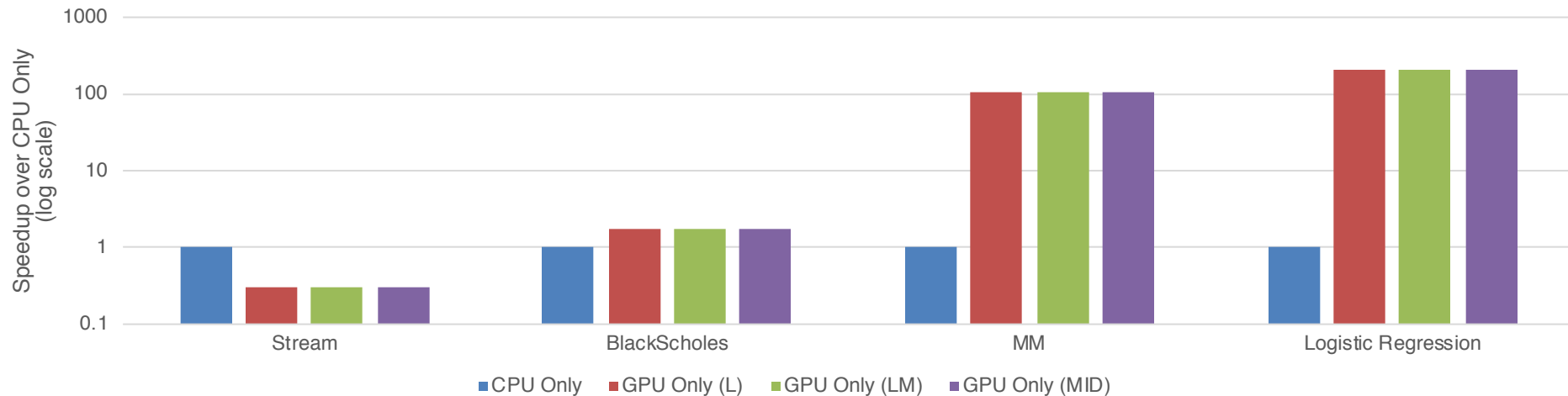
(Multi-nodes, 6GPUs/node,
Cori, Summit)

BlackScholes, MID-level, $n = 2^{30}$



How fast are GPUs?

(A single-node,
Ryzen9 + Radeon RX570)



☐ The GPUAPI works on an AMD GPU machine



A multi-resolution GPU programming model for Chapel

CONCLUSIONS & FUTURE WORK

Conclusions

□ Summary

- The GPUAPI provides an appropriate interface between Chapel and accelerator programs
 - ✓ Source code is available (the feature/explicit branch):
 - <https://github.com/ahayashi/chapel-gpu>
- The LOW-MID and MID level GPU API enable a higher-level GPU programming interface w/ minimal performance overhead
 - ✓ Verified on NVIDIA and AMD GPUs
- The GPUAPI + the GPUIterator module facilitate distributed GPU programming using Chapel
 - ✓ Verified on Summit and Cori
 - ✓ The use of GPUs can significantly improve the performance of Chapel programs



Future Work

□ Building a higher-level programming model

- Built on top of the MID-level API
- forall + “intents”

```
forall i in D with (in: B, out: A) {...}
```

=>

```
var dA = GPUArray(A); var dB = GPUArray(B); toDevice(B);  
kernel(); fromDevice(A);
```

- For the kernel code generation, will explore the possibility of avoiding/minimizing compiler modifications
 - ✓ e.g., generate a C/C++ loop + OpenMP/OpenACC/Other pragma
 - Chapel’s Vectorizing Iterator does a similar thing (vectorizeOnly())

□ Asynchronous GPU API + Futures



Future work (Cont'd)

❑ Wish List: lambda + capture by reference

Current MID-level

```
1 use GPUAPI;
2 var A: [1..n] real(32);
3 var B: [1..n] real(32);
4 var dA = new GPUArray(A);
5 var dB = new GPUArray(B);
6 toDevice(dB);
7 LaunchVC(dA.dPtr(),
           dB.dPtr(), N: size_t);
8 DeviceSynchronize();
9 FromDevice(dA);
10 Free(dA, dB);
```

If dA, dB, N can be captured...

```
1 use GPUAPI;
2 var A: [1..n] real(32);
3 var B: [1..n] real(32);
4 var dA = new GPUArray(A);
5 var dB = new GPUArray(B);
6 var in = (dB); // tuple
7 var out = (dA); // tuple
8 Launch (in, out, lambda () {
9     LaunchVC(dA.dPtr(),
10             dB.dPtr(), N: size_t);
11 });
```



Thank you for your attention!
Any questions?



Backup Slides



Chapel's iterator

- ❑ Chapel's iterator allows us to control over the scheduling of the loops in a productive manner

```
1 // Iterator over fibonacci numbers
2 forall i in fib(10) {
3     A(i) = B(i);
4 }
```

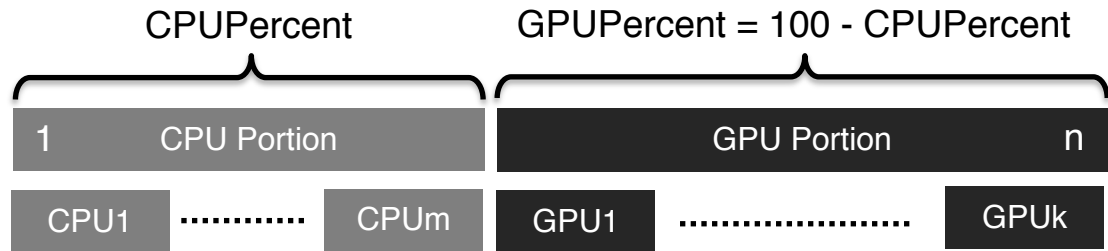
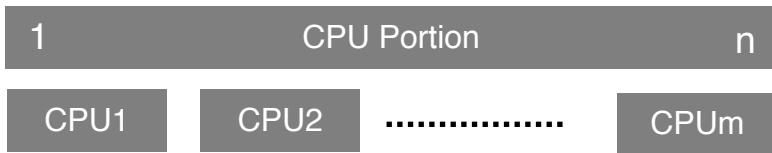
CPU1					CPU2				
0	1	1	2	3	5	8	13	21	34



The GPUterator automates work distribution across CPUs+GPUs

```
1 forall i in 1..n {  
2   A(i) = B(i);  
3 }
```

```
1 forall i in GPU(1..n, GPUWrapper,  
2   CPUPercent) {  
3   A(i) = B(i);  
4 }
```



How to use the GPUIterator?

```
1 var GPUCallback = lambda (lo: int,
2                          hi: int,
3                          nElems: int){
4     assert(hi-lo+1 == nElems);
5     GPUVC(A, B, lo, hi);
6 };
7 forall i in GPU(1..n, GPUCallback,
8                CPUPercent) {
9     A(i) = B(i);
10 }
```

This callback function is called after the GPUIterator has computed the subspace (lo/hi: lower/upper bound, n: # of elements)

GPU() internally divides the original iteration space for CPUs and GPUs



The GPUlterator supports Zippered-forall

```
1 forall (_, a, b) in zip(GPU(1..n, ...), A, B) {  
2     a = b;  
3 }
```

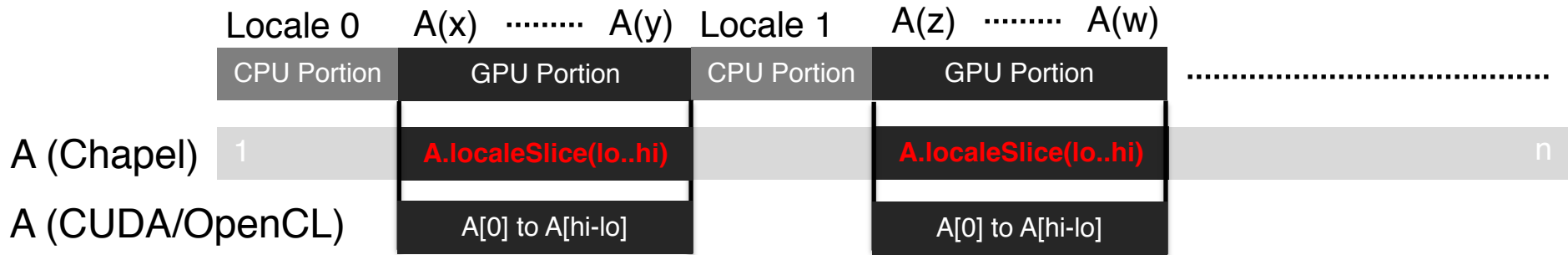
❑ Restriction

- The GPUlterator must be the leader iterator



The GPUerator supports Distributed Arrays (Cont'd)

- ❑ No additional modifications for supporting multi-locale executions



Note: localeSlice is Chapel's array API



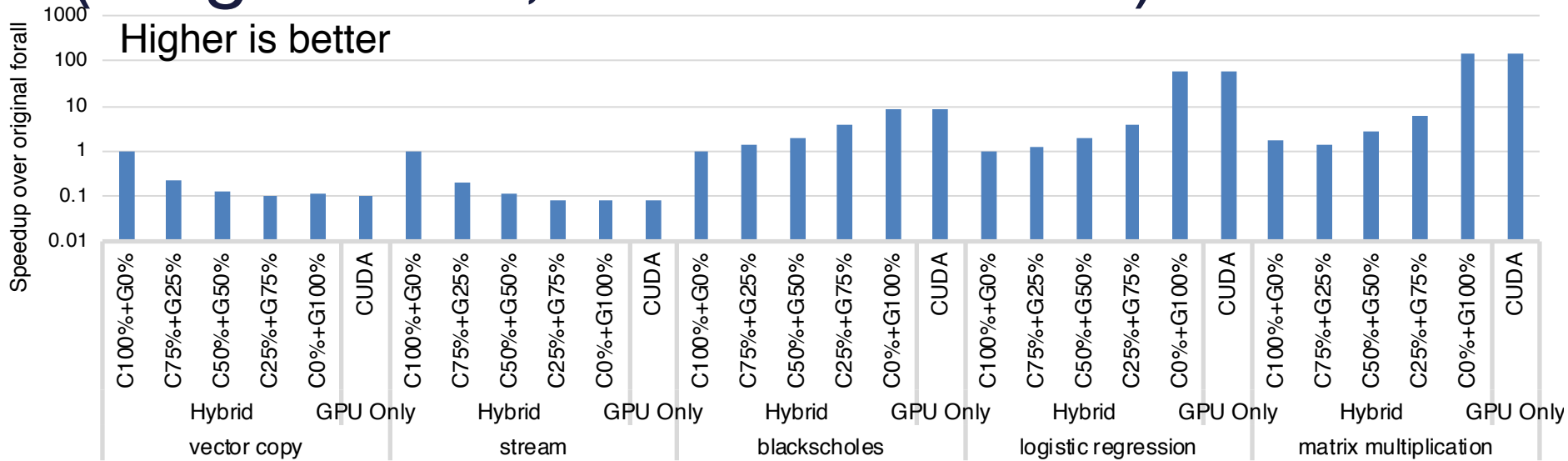
Implementation of the GPUIterator

```
1  coforall subloc in 0..1 {
2    if (subloc == 0) {
3      const numTasks = here.getChild(0).maxTaskPar;
4      coforall tid in 0..#numTasks {
5        const myIters = computeChunk(...);
6        for i in myIters do
7          yield i;
8        }
9      } else if (subloc == 1) {
10     GPUCallback(...);
11   }
12 }
```



How fast are GPUs?

(Single-node, POWER8 + K80)

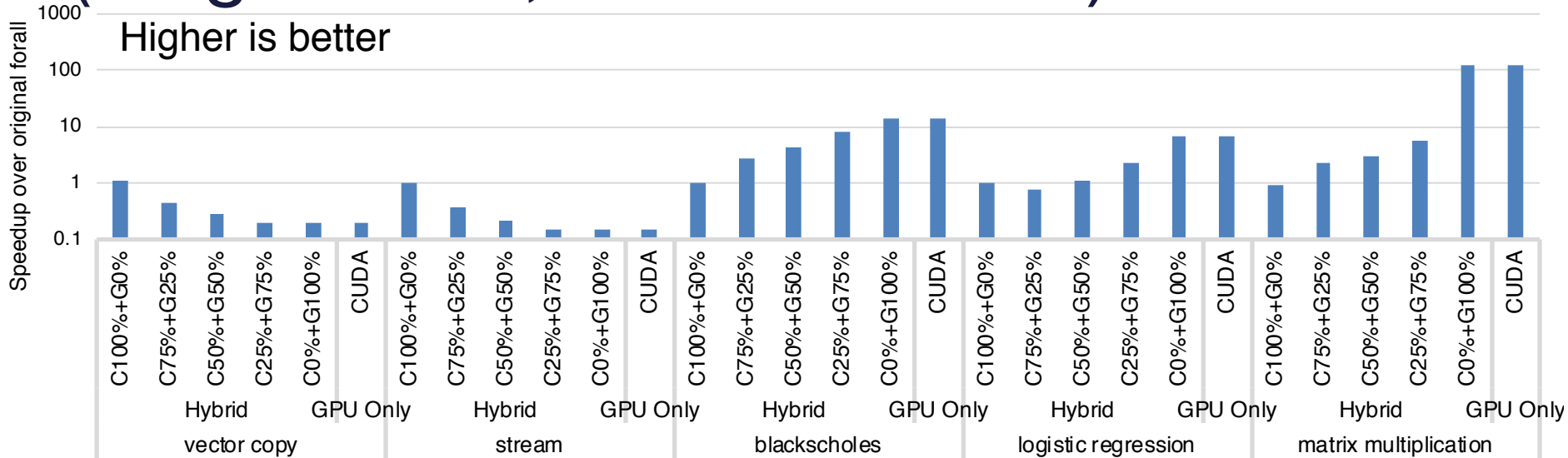


- ❑ The iterator enables exploring different CPU+GPU strategies with very low overheads
- ❑ The GPU is up to 145x faster than the CPU, but is slower than the GPU due to data transfer costs in some cases



How fast are GPUs?

(Single-node, Xeon + M2050)



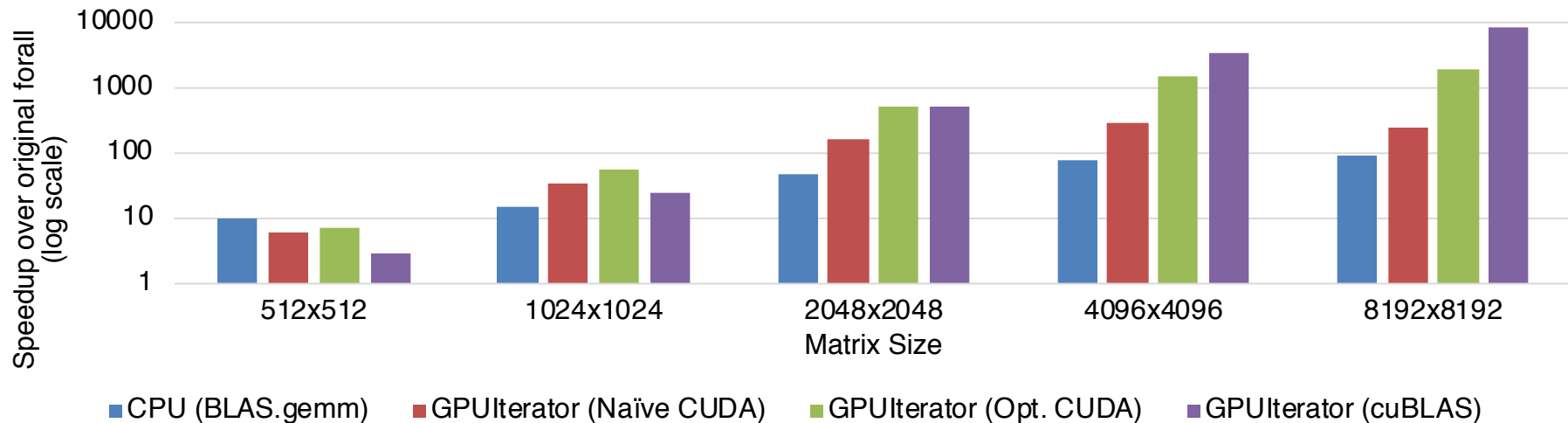
- ❑ The iterator enables exploring different CPU+GPU strategies with very low overheads
- ❑ The GPU is up to 126x faster than the CPU, but is slower than the GPU due to data transfer costs in some cases



How fast are GPUs compared to Chapel's BLAS module on CPUs?

(Single-node, Core i5 + Titan Xp)

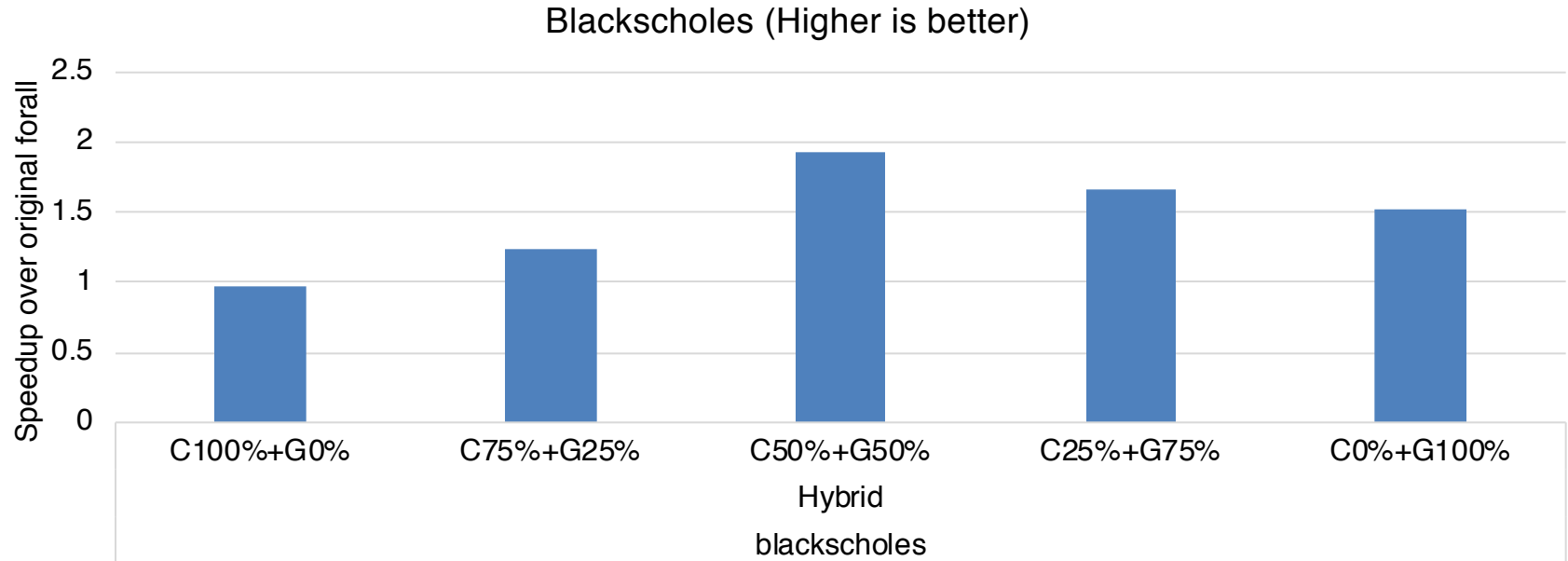
Matrix Multiplication (Higher is better)



- ❑ Motivation: to verify how fast the GPU variants are compared to a highly-tuned Chapel-CPU variant
- ❑ Result: the GPU variants are mostly faster than OpenBLAS's gemm (4 core CPUs)



When is hybrid execution beneficial? (Single node, Core i7+UHD)



- With tightly-coupled GPUs, hybrid execution is more beneficial

