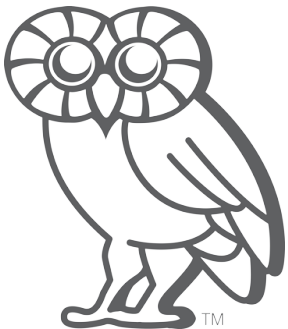


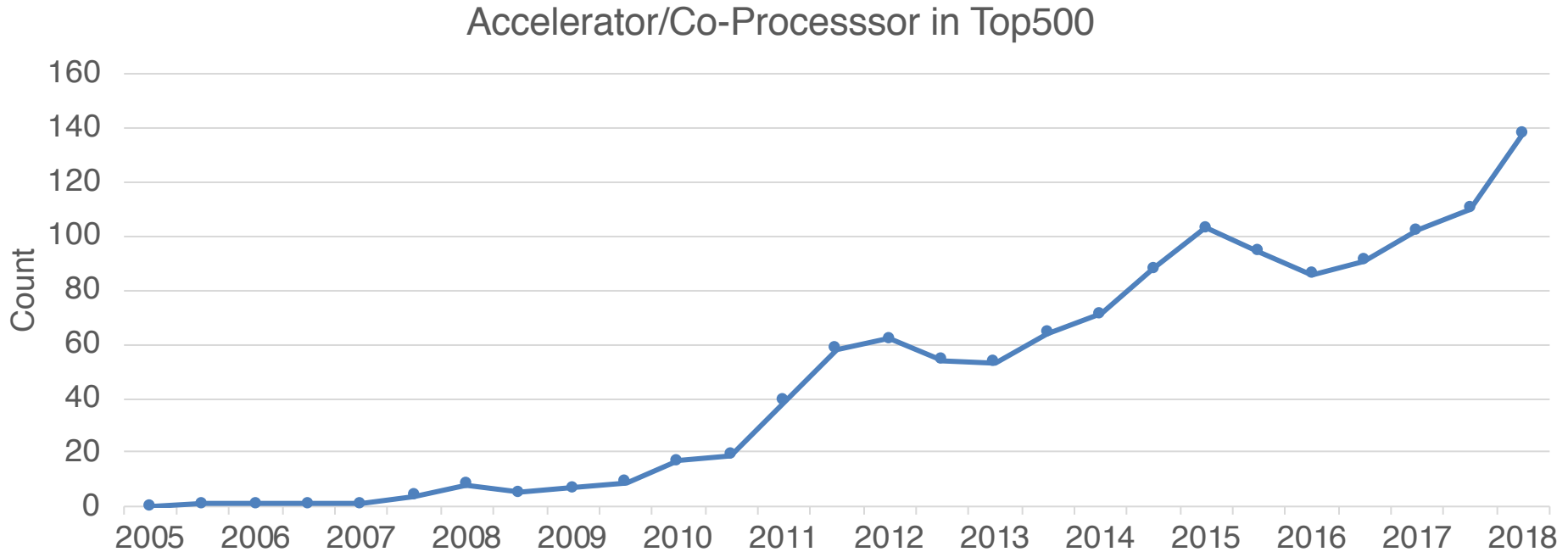


# GPUIterator: Bridging the Gap between Chapel and GPU platforms

Akihiro Hayashi (Rice),  
Sri Raj Paul (Georgia Tech),  
Vivek Sarkar (Georgia Tech)



# GPUs are a common source of performance improvement in HPC



# 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 writing GPU kernels using CUDA/OpenCL or other accelerator language, and invoke them from Chapel (Focus of this paper)



RICE

Georgia  
Tech



[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 (CHI'17).



# Motivation: Vector Copy (Original)

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



# Motivation:

## Vector Copy (GPU)

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

```
1 extern proc GPUVC(A: [] real(32),
2                 B: [] real(32),
3                 lo: int, hi: int);
4
5 var A: [1..n] real(32);
6 var B: [1..n] real(32);
7
8 // Invoking CUDA/OpenCL program
9 GPUVC(A, B, 1, n);
```

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



# Motivation:

## The code is not very portable

```
1 // Original
2 forall i in 1..n {
3     A(i) = B(i);
4 }
```



```
1 // GPU Version
2 GPUVC(A, B, 1, n);
```

### □ Potential “portability” problems

- How to switch back and forth between the the original version and the GPU version?
- How to support hybrid execution?
- How to support distributed arrays?

### Research Question:

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



# Our Solution: GPUIterator

```
1 // Original Version
2 forall i in 1..n {
3     A(i) = B(i);
4 }
```



```
1 // GPU Version
2 GPUVC(A, B, 1, n);
```

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

## □ Contributions:

- Design and implementation of the GPUIterator
- Performance evaluation of different CPU+GPU execution strategies



# 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

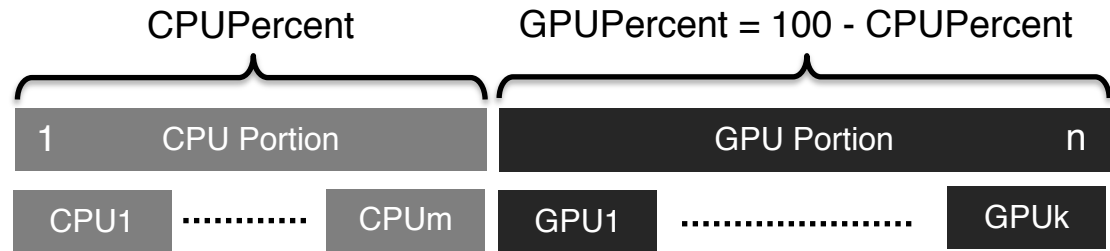
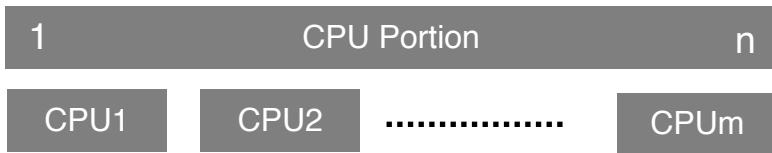




# The GPUterator automates work distribution across CPUs+GPUs

```
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 GPUterator supports Distributed Arrays

```
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 GPUCallback = lambda (lo: int, hi: int, nElems: int) {
5      GPUVC(A.localSlice(lo..hi),
6            B.localSlice(lo..hi),
7            0, hi-lo, nElems);
8  };
9  forall i in GPU(D, GPUCallback,
10                CPUPercent) {
11      A(i) = B(i);
12  }
```



# The GPUlterator supports Zippered-forall

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

## ❑ Restriction

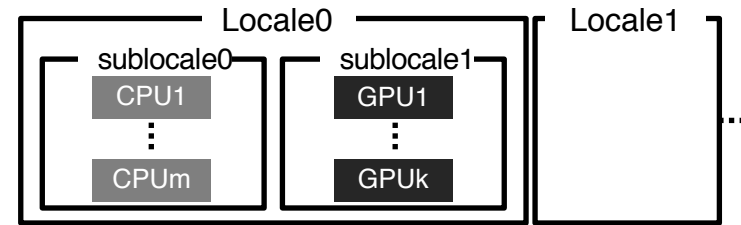
- The GPUlterator must be the leader iterator



# Implementation of the GPUIterator

## ❑ Internal modules

- <https://github.com/ahayashi/chapel>
- Created the GPU Locale model
  - ✓ CHPL\_LOCALE\_MODEL=gpu



## ❑ External modules

- <https://github.com/ahayashi/chapel-gpu>
- Fully implemented in Chapel



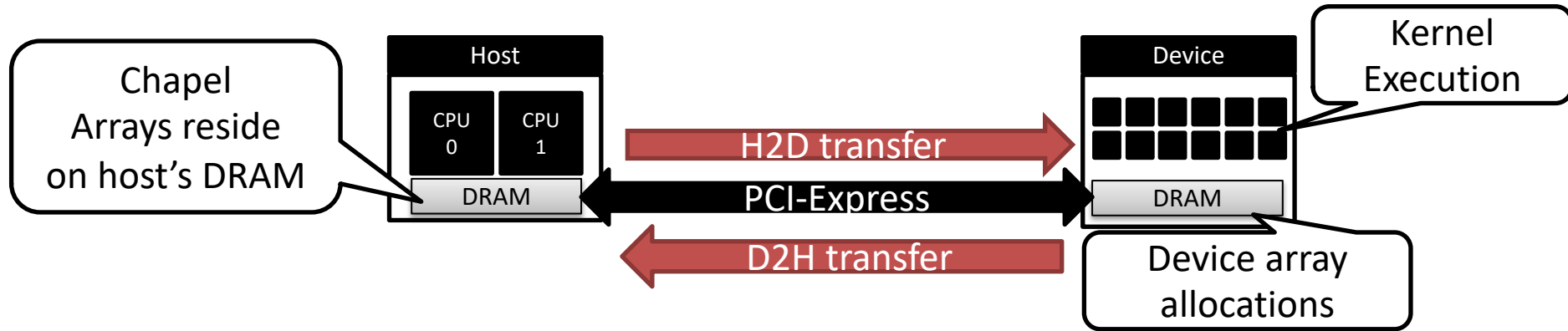
# Implementation of the GPUIterator

```
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 }
```



# Writing CUDA/OpenCL Code for the GPUterator

- GPU programs for the GPUterator should include typical host and device operations



# Performance Evaluations

## ❑ Platforms

- Intel Xeon CPU (12 cores) + NVIDIA Tesla M2050 GPU
- IBM POWER8 CPU (24 cores) + NVIDIA Tesla K80 GPU
- Intel Core i7 CPU (6 cores) + Intel UHD Graphics 630/AMD Radeon Pro 560X
- Intel Core i5 CPU (4 cores) + NVIDIA TITAN Xp

## ❑ Chapel Compilers & Options

- Chapel Compiler 1.20.0-pre (as of March 27) with the --fast option

## ❑ GPU Compilers

- CUDA: NVCC 7.0.27(M2050), 8.0.61 (K80) with the -O3 option
- OpenCL: Apple LLVM 10.0.0 with the -O3 option





# Performance Evaluations (Cont'd)

- ❑ Tasking
  - CUDA: CHPL\_TASK=qthreads
  - OpenCL: CHPL\_TASK=fifo
- ❑ Applications (<https://github.com/ahayashi/chapel-gpu>)
  - Vector Copy
  - Stream
  - BlackScholes
  - Logistic Regression
  - Matrix Multiplicaiton



# How many lines are added/modified?

	LOC added/modified (Chapel)	CUDA LOC (for NVIDIA GPUs)	OpenCL LOC (for Intel/AMD GPUs)
Vector Copy	6	53	256
Stream	6	56	280
BlackScholes	6	131	352
Logistic Regression	11	97	472
Matrix Multiplication	6	213	290

Source code changes are minimal



RICE

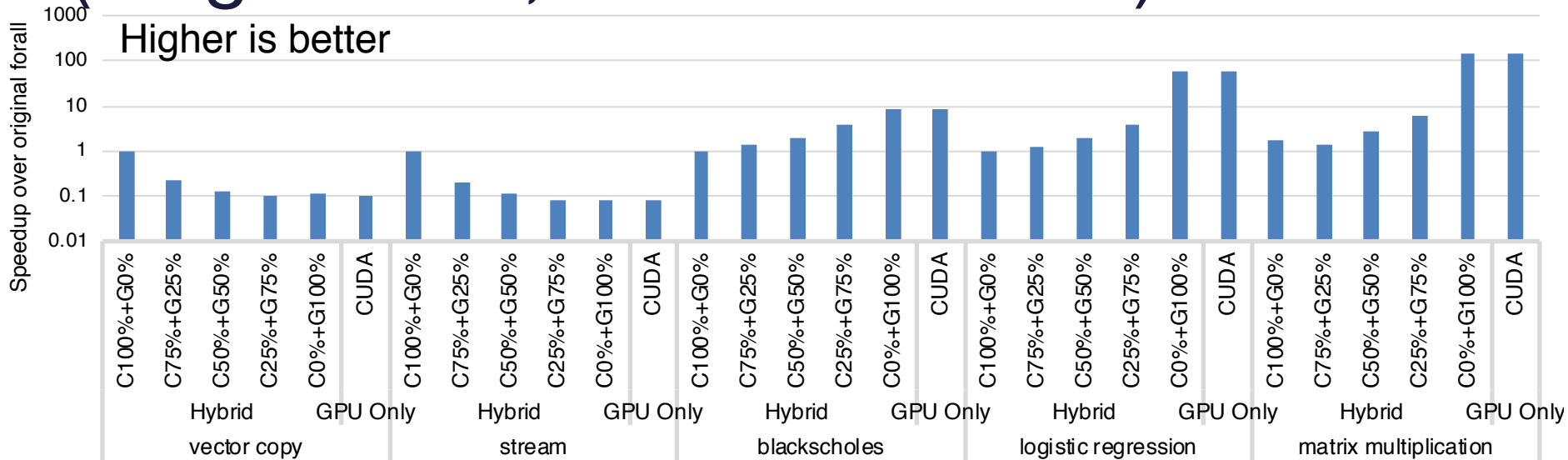


LOC for CUDA/OpenCL are out of focus



# How fast are GPUs?

## (Single-node, POWER8 + K80)

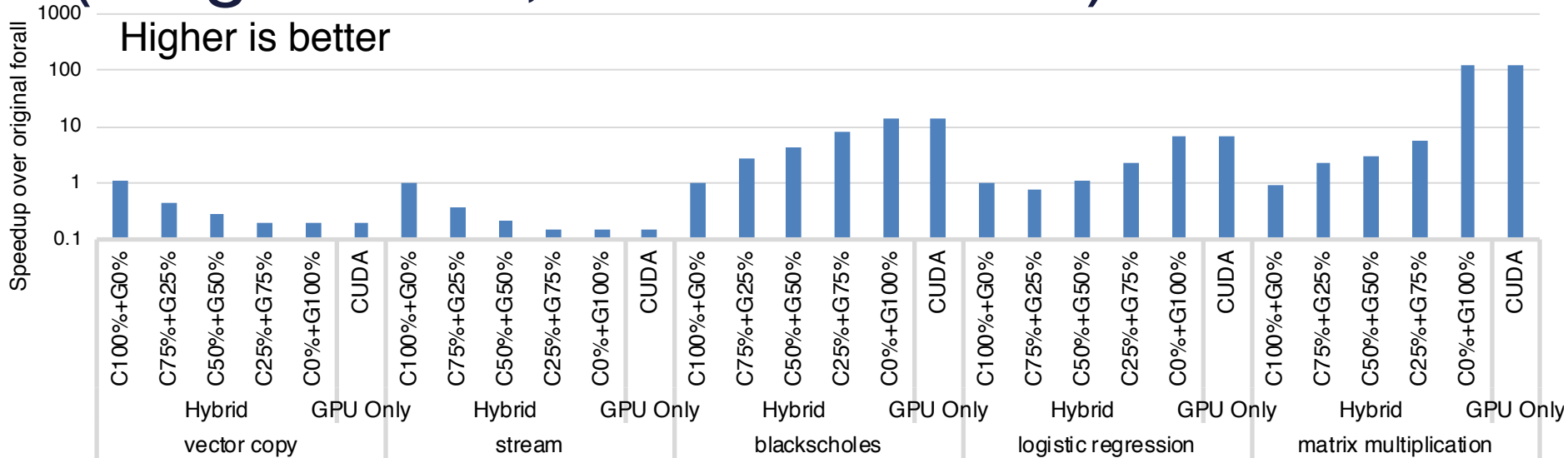


- ❑ The iterator enables exploring different CPU+GPU strategies with very low overheads
- ❑ The GPU is up to 145x faster than the CPU, but is slower than the GPU due to data transfer costs in some cases



# How fast are GPUs?

## (Single-node, Xeon + M2050)



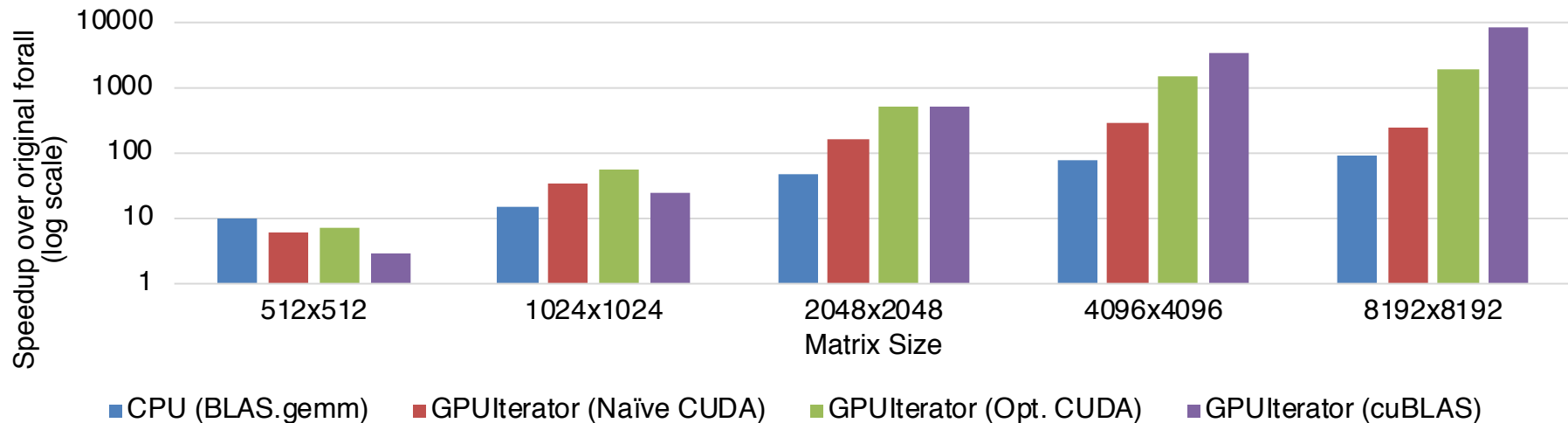
- ❑ The iterator enables exploring different CPU+GPU strategies with very low overheads
- ❑ The GPU is up to 126x faster than the CPU, but is slower than the GPU due to data transfer costs in some cases



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

(Single-node, Core i5 + Titan Xp)

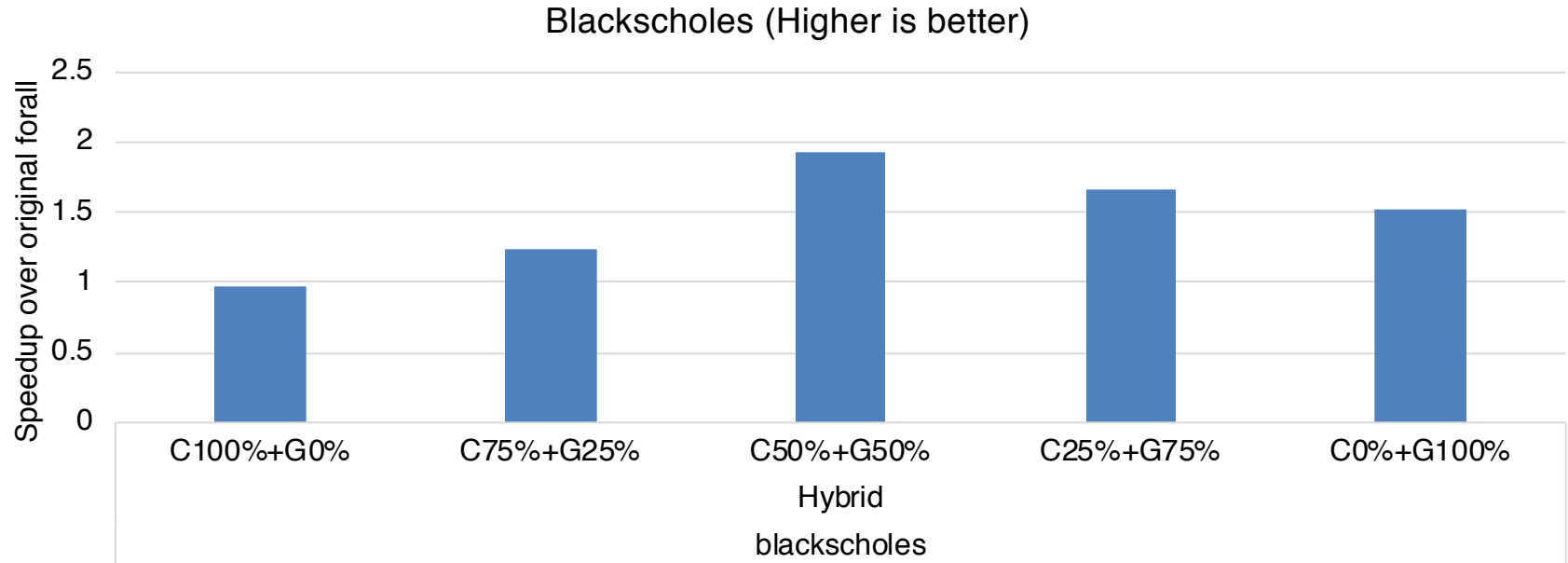
Matrix Multiplication (Higher is better)



- ❑ 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)



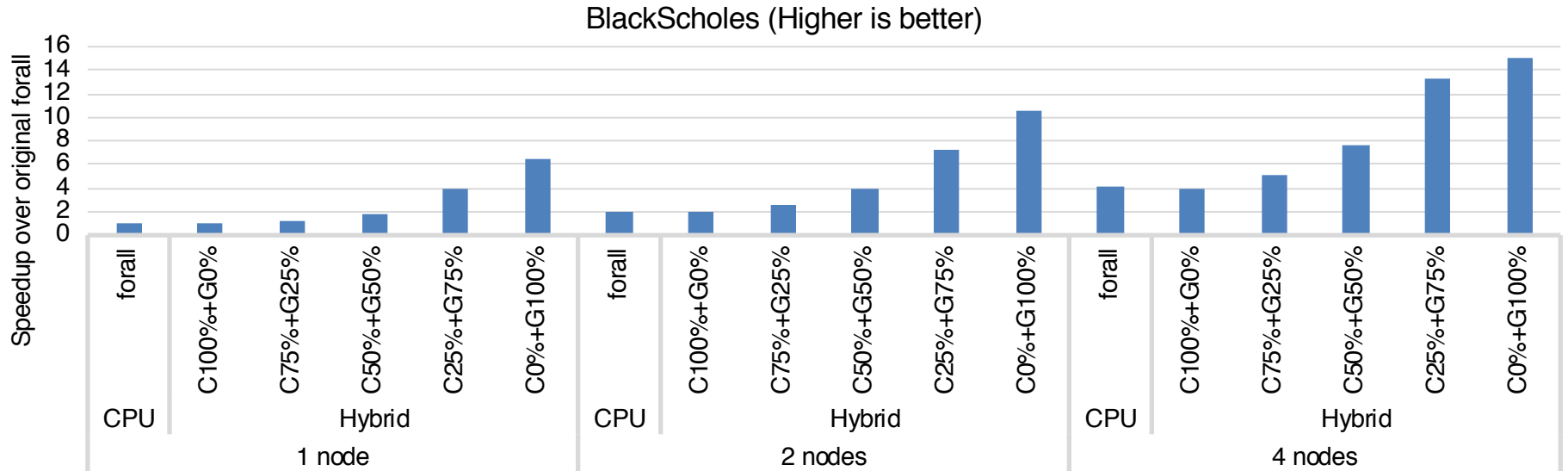
# When is hybrid execution beneficial? (Single node, Core i7+UHD)



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



# Multi-node performance numbers (Xeon + M2050)



- ❑ The original forall show good scalability
- ❑ The GPU variants give further performance improvements



# Conclusions & Future Work

## □ Summary

- The GPUlterator provides an appropriate interface between Chapel and accelerator programs
  - ✓ Source code is available:
    - <https://github.com/ahayashi/chapel-gpu>
- The use of GPUs can significantly improve the performance of Chapel programs

## □ Future Work

- Support reduction
- Further performance evaluations on multi-node CPU+GPU systems
- Automatic selection of the best “CPUPercent”

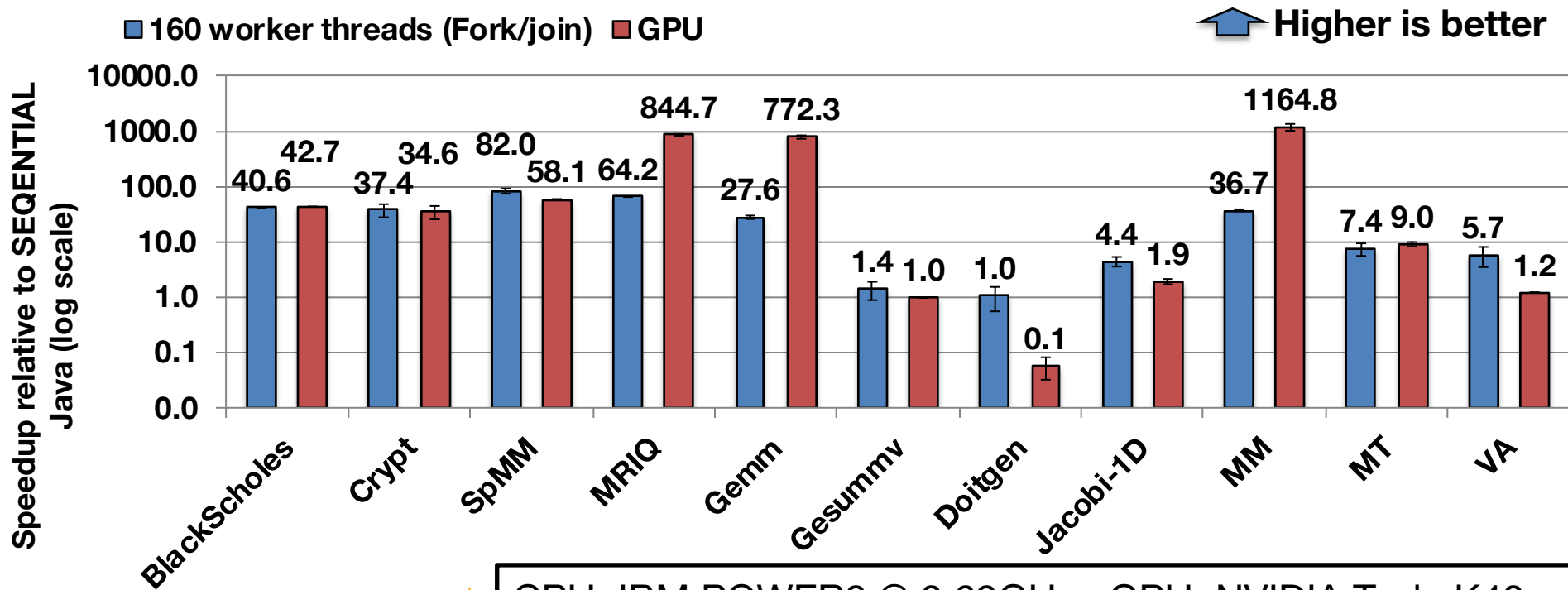




# Backup Slides



# GPU is not always faster



RICE

Georgia Tech



CPU: IBM POWER8 @ 3.69GHz , GPU: NVIDIA Tesla K40m

# The GPUerator supports Distributed Arrays (Cont'd)

- ❑ No additional modifications for supporting multi-locale executions

