

# Enabling CHIP-SPV in Chapel GPUAPI module

CHI UW2023, June 2, 2023

Jisheng Zhao (Georgia Tech)

Akihiro Hayashi (Georgia Tech)

Brice Videau (ANL)

Vivek Sarkar (Georgia Tech)



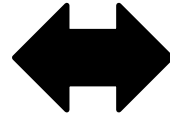
Enabling CHIP-SPV in Chapel GPUAPI module

# MOTIVATION

# GPU Programming in Chapel: no “intermediate” programming model

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

HIGH-level:  
The compiler compiles  
forall to CUDA, HIP, and  
OpenCL

forall

The missing link

LOW-level:  
The user prepares full  
GPU programs and  
invokes them from Chapel  
(w/ or w/o the GPUerator)

C GPUerator [1]  
Interoperability  
CUDA/HIP/OpenCL  
NVIDIA/AMD/Other GPUs

Our proposal

**Chapel programmer friendly GPU APIs :**  
MID-level  
`var dA = new GPUArray(A);`  
`dA.toDevice();`

**Thin wrappers for low-level GPU APIs:**  
MID-LOW-level  
`Malloc(); Memcpy();`

Goal: increase productivity with no performance loss



# Contributions

## ❑ Why higher-level abstraction of GPU API?

- For improving productivity

Our observation: The complexity in GPU programming comes not only from writing GPU kernels in the device part, but also from writing the host part

- ✓ Our GPUAPI is designed to simplify the host part

- For improving portability

Our observation: There are different GPU programming models from different vendors

- ✓ Our GPUAPI is implemented to work on different platforms (NVIDIA, AMD, Intel, ...)

## ❑ Contribution (specific to this talk):

- Enhance support for Intel GPUs by implementing a CHIP-SPV backend in the GPUAPI module



Enabling CHIP-SPV in Chapel GPUAPI module

# THE GPUAPI MODULE

# Summary of the Chapel GPUAPI module

- ❑ 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
  - MID-LOW: 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);`
- ❑ Note
  - The user is still supposed to write kernels in CUDA/HIP/SYCL (DPC++)
  - 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
  - The module can be utilized in real-world Chapel applications such as Champs and ChOp



# 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 }
```

Chapel's  
Distributed Array Allocation  
(n divided by # of nodes float elements)

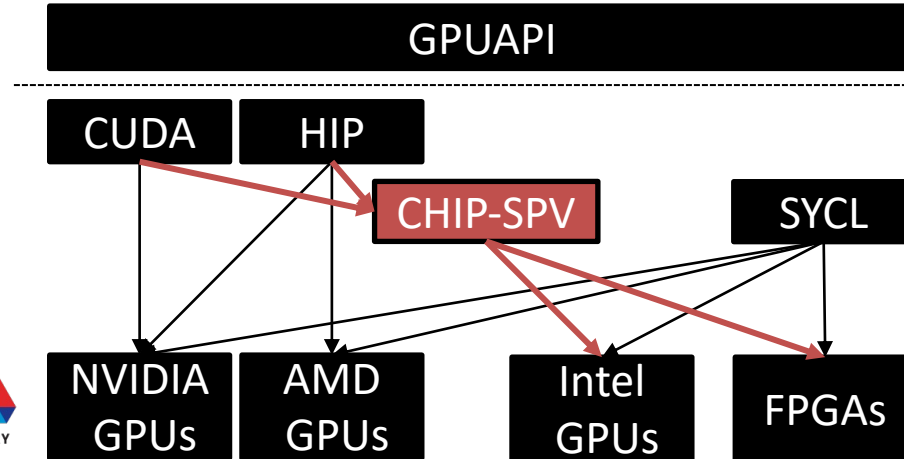
```
1 // separate C file (CUDA w/ device lambda)
2 void LaunchST(float *dA, float *dB,
3              float *dC, float alpha, int N) {
4     GPU_FUNCTOR(N, 1024, NULL,
5                 [=] __device__ (int i) {
6         dA[i] = dB[i] + alpha * dC[i];
7     });
8 }
```

The user has the option of writing device functions, device lambdas, or library calls

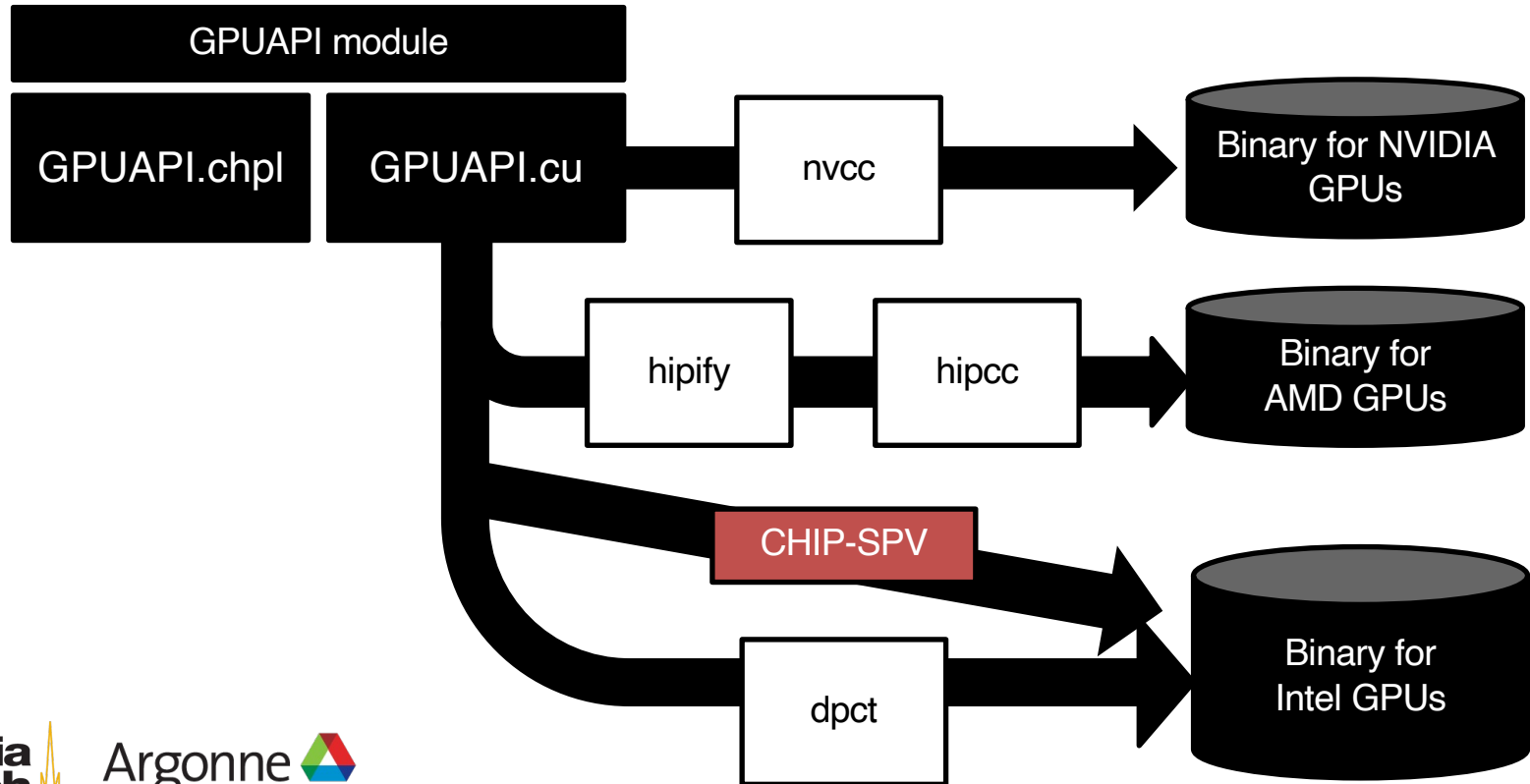


# Multi-platform Support: Module implementation

- ❑ Our module supports a wide variety of GPUs
- ❑ Our cmake-based build system detects types of GPUs and generate corresponding static and shared libraries
  - libGPUAPIX\_static.a and libGPUAPIX.so
    - ✓ X = CUDA, HIP, or DPCPP



# Multi-platform Support: Module implementation



# Multi-platform Support: User-written kernels

- ❑ Our multi-platform GPUAPI support allows the user to choose any GPU programming models to write their kernel code

	CUDA	HIP	SYCL
NVIDIA	Yes	Yes	Yes
AMD	Yes (hipify)	Yes	Yes
Intel	Yes (dpct or CHIP-SPV)	Yes (CHIP-SPV)	Yes

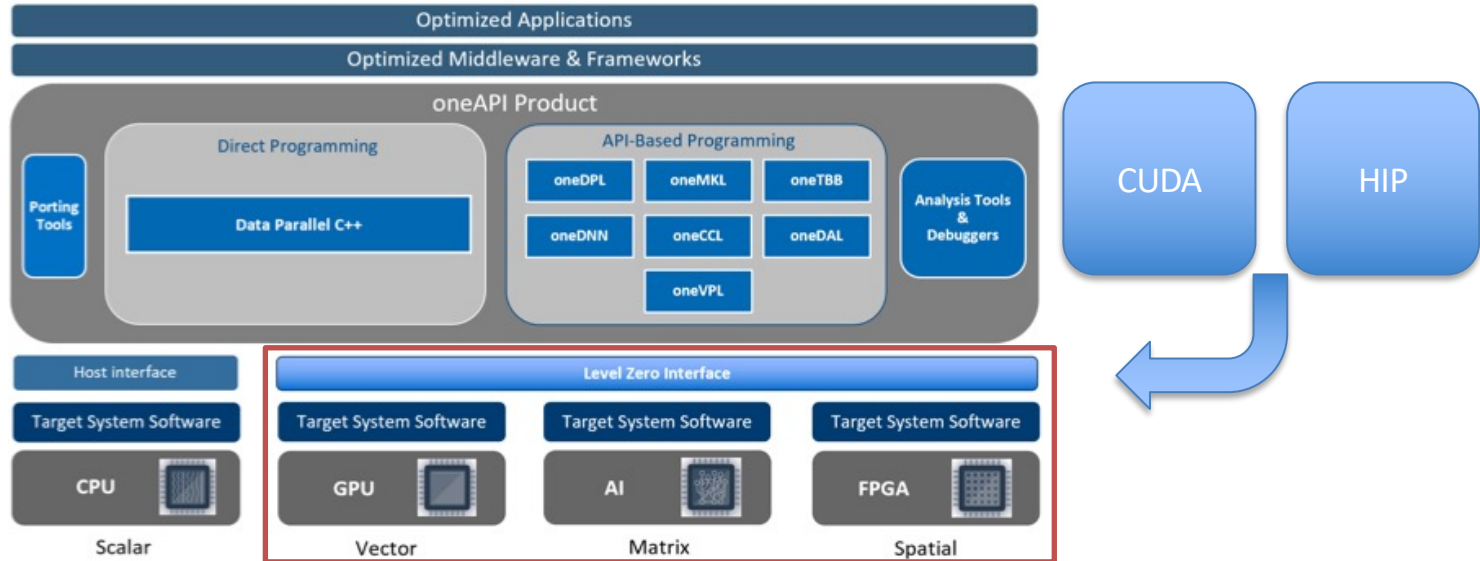


Enabling CHIP-SPV in Chapel GPUAPI module

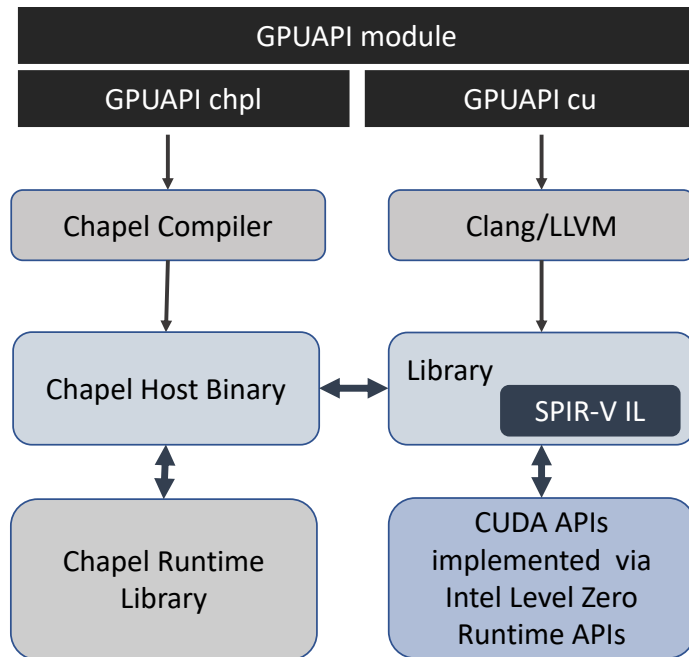
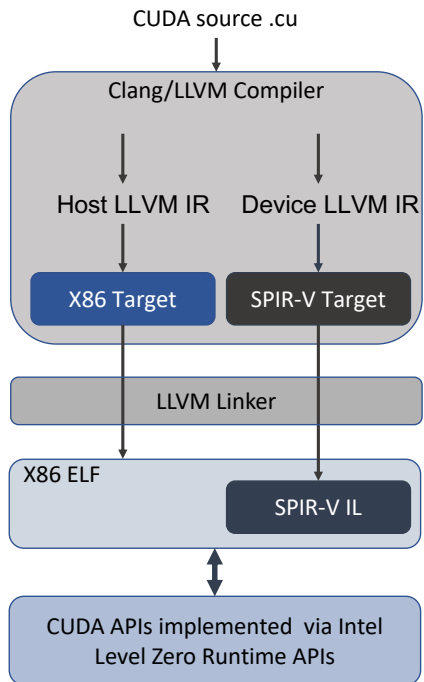
# CHIP-SPV

# CHIP-SPV

- ❑ Connect heterogeneous programming models (CUDA/HIP) to low-level Intel Level Zero runtime
- ❑ Intel Level Zero: A system level programming interface that bridges high level libraries to Intel devices



# CHIP-SPV in the GPUAPI module



Enabling CHIP-SPV in Chapel GPUAPI module

# PRELIMINARY PERFORMANCE EVALUATIONS

# Performance Evaluations

## ❑ Platforms (Single-node, integrated-GPUs)

- Intel Xeon Processor E5-1585 v5 (4-core) + Iris Pro P580 (Gen9)
- Intel i7-12700 + UHD770

## ❑ Applications

- Micro-benchmark: STREAM, BlackScholes, Matrix Multiplication, Logistic Regression

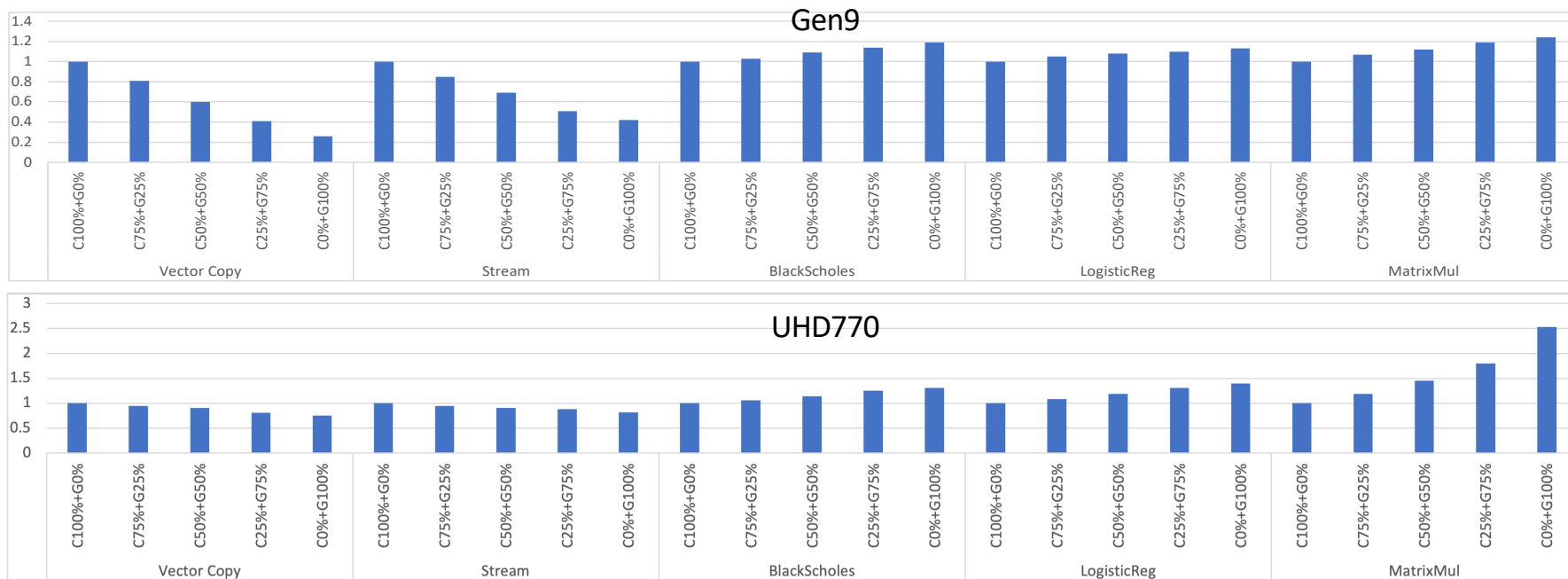
## ❑ Chapel Compilers & Options

- Chapel Compiler 1.29.0 with the --fast option
- Used with the GPUlterator (CPU+GPU execution)





# Preliminary Performance Numbers (vs. forall execution on CPUs)



## □ Key take aways:

- Our end-to-end compilation flow with CHIP-SPV is verified
- Performance improvements are not significant due to integrated GPUs, but this is where the GPUlatorator module comes into play

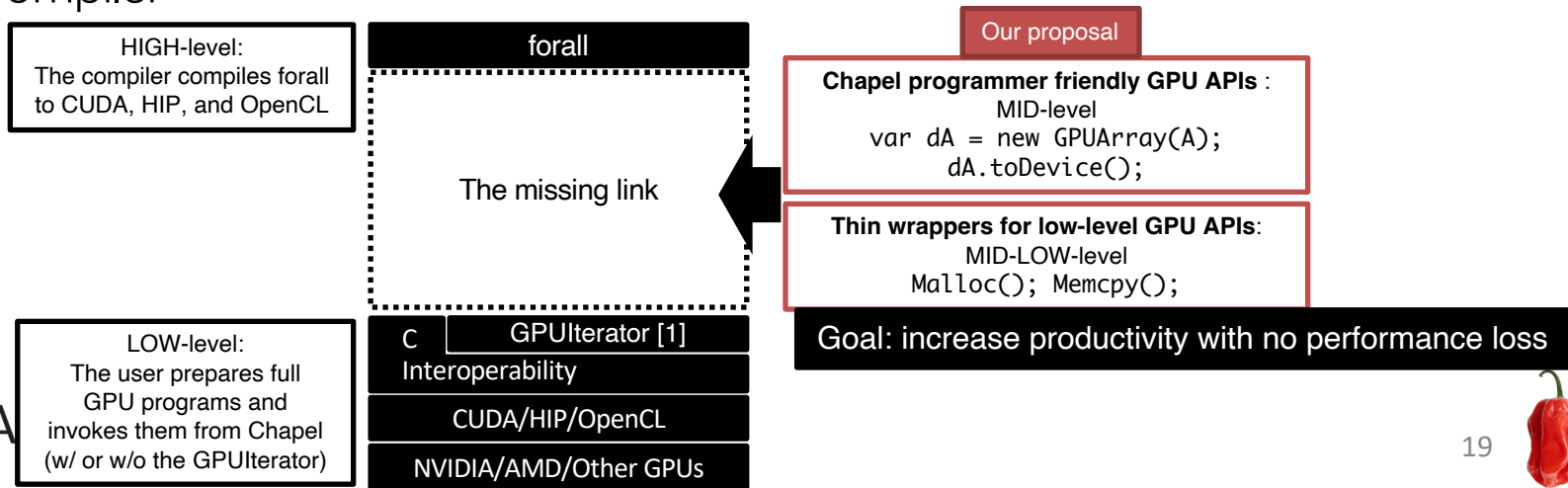


Enabling CHIP-SPV in Chapel GPUAPI module

# CONCLUSIONS

# Conclusions

- ❑ Intel GPUs support in the GPUAPI module
  - Verified with microbenchmarks on two intel GPU platforms
- ❑ Future work:
  - Use discrete Intel GPUs
  - Discuss the possibility of using CHIP-SPV in the current Chapel compiler



# Join our community

## □ GPUAPI+GPUlterator:

- The repository
  - ✓ <https://github.com/ahayashi/chapel-gpu>
- Detailed Documents
  - ✓ <https://ahayashi.github.io/chapel-gpu/index.html>

## □ Our community is growing!



Contact: ahayashi "at" gatech.edu



# Acknowledgements

- This work was supported by the Argonne Leadership Computing Facility, which is a DOE Office of Science User Facility supported under Contract DE-AC02-06CH11357, and by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of two U.S. Department of Energy organizations (Office of Science and the National Nuclear Security Administration). We also gratefully acknowledge the computing resources provided and operated by the Joint Laboratory for System Evaluation (JLSE) at Argonne National Laboratory.



Enabling CHIP-SPV in Chapel GPUAPI module

# BACKUP SLIDES

# Chapel GPU API Design:

## MID-LOW 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(...)`;
  - ✓ `MallocPitch(...)`;
- Host-to-device, and device-to-host data transfers
  - ✓ `Memcpy(...)`;
  - ✓ `Memcpy2D(...)`;
- Ensuring the completion of GPU computations
  - ✓ `DeviceSynchronize()`;
- Device Memory deallocation
  - ✓ `Free(...)`;



# 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);`
  - ✓ `var dA = new GPUJaggedArray(A);`
- Host-to-device, and device-to-host data transfers
  - ✓ `ToDevice(dA:GPUArray, ...); FromDevice(dA: GPUArray, ...);`
  - ✓ `dA.ToDevice(); dA.fromDevice();`
- Implicit Device Memory deallocation
  - ✓ Automatically “freed” when a GPUArray/GPUJaggedArray object is deleted
- Explicit Device Memory deallocation
  - ✓ `delete`





# Chapel GPU API Design:

## MID-LOW/MID GPU API Example

### MID-LOW 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, TODEVICE);
12 Memcpy(dC, c_ptrTo(C), size, TODEVICE);
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);
```

Chapel's  
Array Allocation  
(n float elements)

### 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);
```

Chapel's  
Array Allocation  
(n float elements)



# Our GPUerator module facilitates running GPUAPI programs across CPUs+GPUs

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

```
1 forall i in GPU(1..n, GPUCallback,  
2   CPUPercent) {  
3   A(i) = B(i) + alpha * C(i);  
4 }
```

GPUCallback is a function that includes a sequence of GPUAPI (next slide)

