

# Recent GPU Improvements in Charm++ 6.9

Jæmin Choi, **Michæel Robson**

Parallel Programming Laboratory  
University of Illinois Urbana-Champaign

November 14, 2018

## How to Utilize GPUs in Charm++

1. Use CUDA directly
  - ▶ Let each chare offload (small) kernels
  - ▶ Or manually aggregate data at a synchronization point and offload one big kernel
2. Use **GPU Manager** library of Charm++
  - ▶ Why? What good is it?

## Problems with Using GPUs in Charm++

- ▶ Due to **overdecomposition** and **asynchrony**
  1. Granularity of work
  2. Blocking offload API
  3. Responsiveness

## Current GPU Manager

- ▶ Addresses Problem 2 (blocking offload API)
- ▶ User constructs and submits a `WorkRequest` object, specifying
  - ▶ Data buffers and directions of transfer
  - ▶ Kernel to be executed and its specifications (e.g. grid size, block size)
- ▶ Runtime tracks `WorkRequests`, overlapping data transfers with kernel execution
  - ▶ But does NOT overlap multiple kernel executions
  - ▶ Because only one CUDA stream is used for kernels
- ▶ Execution continues without blocking after `WorkRequest` submission
- ▶ 3 CUDA streams used internally: Data-in, Kernel, Data-out
- ▶ **Problems**
  - ▶ Only one CUDA stream for all kernels
  - ▶ Unnecessarily complex API

## New GPU Manager: Release 6.9.0

- ▶ Partially addresses Problem 1 (granularity of work)
  - ▶ Allows kernels to execute in separate CUDA streams
  - ▶ Runtime support for kernel aggregation is ongoing research
- ▶ Non-blocking feature implemented using CUDA events
- ▶ Much simpler API (almost identical to CUDA API)
  - ▶ **Hybrid API:** `hapi` prefix instead of `cuda`
  - ▶ `hapiAddCallback()`: invoke provided Charm++ callback function when data transfer/kernel execution completes, replaces `cudaStreamSynchronize()`
- ▶ Ongoing research to address Problem 3 (responsiveness)

## Non-blocking Implementation of Offloading

- ▶ Use CUDA events to detect completion of GPU work
- ▶ Each PE maintains a queue of events
- ▶ Queue is checked in the scheduler before choosing what to execute next
- ▶ Charm++ callback invoked on completion to continue program flow
- ▶ Impractical for the user to implement
  - ▶ Unclear where in the program flow the queue should be checked
  - ▶ Unclear how frequent the checking should occur
- ▶ Alternative: CUDA callback, but single callback thread becomes a bottleneck

# Matmul Code Comparison: CUDA, New GPU Manager

```
2 void cudaMatMul(ElementType *h_A, ElementType *h_B, ElementType *h_C,  
3               ElementType *d_A, ElementType *d_B, ElementType *d_C,  
4               cudaStream_t stream, int matrixSize) {  
5     int size = matrixSize * matrixSize * sizeof(ElementType);  
6     dim3 block(BLOCK_SIZE, BLOCK_SIZE);  
7     dim3 grid(ceil((float)matrixSize / block.x),  
8             ceil((float)matrixSize / block.y));  
9  
10    cudaMemcpyAsync(d_A, h_A, size, cudaMemcpyHostToDevice, stream);  
11    cudaMemcpyAsync(d_B, h_B, size, cudaMemcpyHostToDevice, stream);  
12  
13    matrixMul<<<grid, block, 0, stream>>>(d_C, d_A, d_B, matrixSize, matrixSize);  
14  
15    cudaMemcpyAsync(h_C, d_C, size, cudaMemcpyDeviceToHost, stream);  
16  
17    cudaStreamSynchronize(stream);  
18 }
```

Figure: CUDA

```
2 void cudaMatMul(ElementType *h_A, ElementType *h_B, ElementType *h_C,  
3               ElementType *d_A, ElementType *d_B, ElementType *d_C,  
4               void *cb, int matrixSize) {  
5     int size = matrixSize * matrixSize * sizeof(ElementType);  
6     dim3 block(BLOCK_SIZE, BLOCK_SIZE);  
7     dim3 grid(ceil((float)matrixSize / block.x),  
8             ceil((float)matrixSize / block.y));  
9  
10    cudaStream_t stream = hapiGetStream();  
11  
12    hapiCheck(hapiMemcpyAsync(d_A, h_A, size, cudaMemcpyHostToDevice, stream));  
13    hapiCheck(hapiMemcpyAsync(d_B, h_B, size, cudaMemcpyHostToDevice, stream));  
14  
15    matrixMul<<<grid, block, 0, stream>>>(d_C, d_A, d_B, matrixSize, matrixSize);  
16    hapiCheck(cudaPeekAtLastError());  
17  
18    hapiCheck(hapiMemcpyAsync(h_C, d_C, size, cudaMemcpyDeviceToHost, stream));  
19  
20    hapiAddCallback(stream, cb);  
21 }
```

Figure: New API

## Performance Evaluation: busywait

- ▶ Benchmark designed to validate new GPU Manager
- ▶ Tasks (kernels on GPU) busywait both on CPU and GPU
- ▶ Vary how much work out of total is offloaded, and how long they take
- ▶ 3 configurations of task duration:
  - ▶ CPU 1 ms, GPU 10 ms
  - ▶ CPU 10 ms, GPU 1 ms
  - ▶ CPU 10 ms, GPU 10 ms
- ▶ 8 PEs, 16 chares per PE, 128 chares total, 100 iterations
- ▶ 32 concurrent kernels with new GPU Manager (vs. 8 without)
- ▶ Up to **4.79x** speedup compared to directly using CUDA
- ▶ Effectiveness of runtime support depends on application characteristics



## Performance Evaluation: busywait

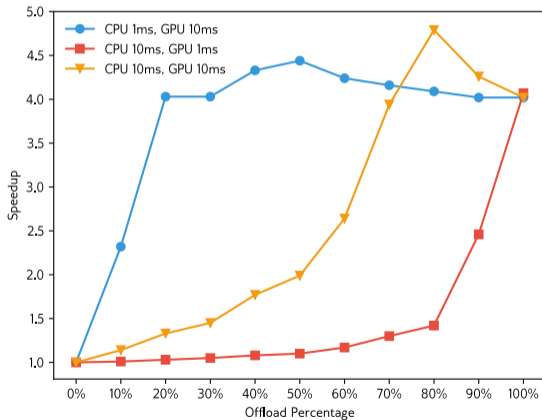


Figure: Speedup of busywait benchmark

## Performance Evaluation: stencil2d

- ▶ 2D 5-point iterative stencil benchmark
- ▶ Evaluate effectiveness under realistic workload
- ▶ 16,384 x 16,384 grid, decomposed into 512 x 512 blocks (chares)
- ▶ 8 PEs, 128 chares per PE, 1,024 chares total, 100 iterations
- ▶ Vary percentage of chares that offload work to GPU
- ▶ 32 concurrent kernels with new GPU Manager (vs. 8 without)
- ▶ Up to **2.75x** speedup compared to directly using CUDA

## Performance Evaluation: stencil2d

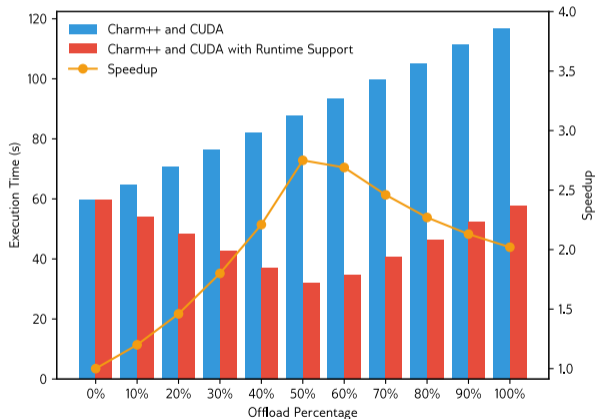


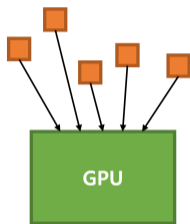
Figure: Execution Time and Speedup of stencil2d benchmark

## Conclusion

- ▶ New GPU Manager: presented as a ACM SRC poster at SC'17
- ▶ 3 main issues with using GPUs in Charm++
  1. Granularity
  2. Blocking
  3. Responsiveness
- ▶ Mostly resolved issue #2, but need more work on issues #1 and #3
- ▶ Interesting research topics with fine-grain tasks and GPUs
- ▶ Increasing importance of accelerators even for irregular applications

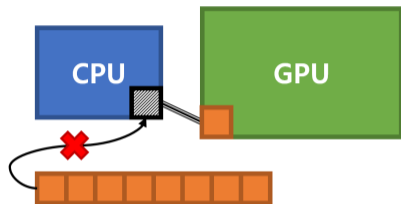
# Backup Slides

## Problem 1: Granularity of Work



- ▶ Each chare is fine-grained
- ▶ Contain little data and work → small kernels
- ▶ Kernels should be able to execute concurrently
- ▶ Or need to aggregate kernels

## Problem 2: Blocking Offload API



- ▶ Commonly used CUDA API are blocking
  - ▶ E.g. `cudaDeviceSynchronize()`, `cudaStreamSynchronize()`
- ▶ PEs are implemented as persistent threads on CPU cores
- ▶ Blocking call thus prevents another chore from executing
- ▶ Another problem: number of concurrent kernels limited to the number of PEs
- ▶ Offload API should be **non-blocking** for Charm++

## Problem 3: Responsiveness

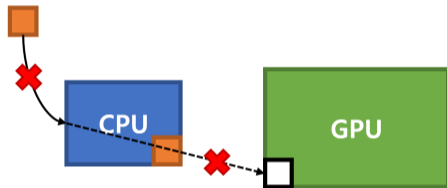


Figure: Slow initiation

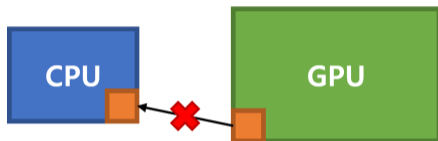


Figure: Slow response

1. Slow initiation
  - ▶ Method offloading work must wait if target PE is busy (even if the GPU is free)
2. Slow response
  - ▶ Handling completed GPU work delayed if target PE is busy



## Performance Evaluation: Test Environment

- ▶ Single compute node of OLCF Titan
- ▶ Up to 8 cores of AMD Opteron 6274 CPU
- ▶ 32GB DDR3 memory
- ▶ NVIDIA Tesla K20X GPU

## GPU Applications: ChaNGa

- ▶ Cosmological N-body simulations
- ▶ Leverages GPU Manager
- ▶ Offloads physics kernels
- ▶ Active work in optimization

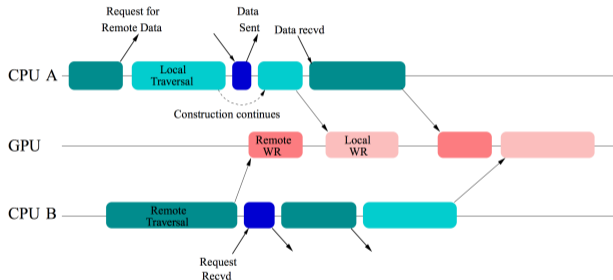


Figure: ChaNGa GPU Manager Design

## GPU Applications: Recent ChaNGa Results

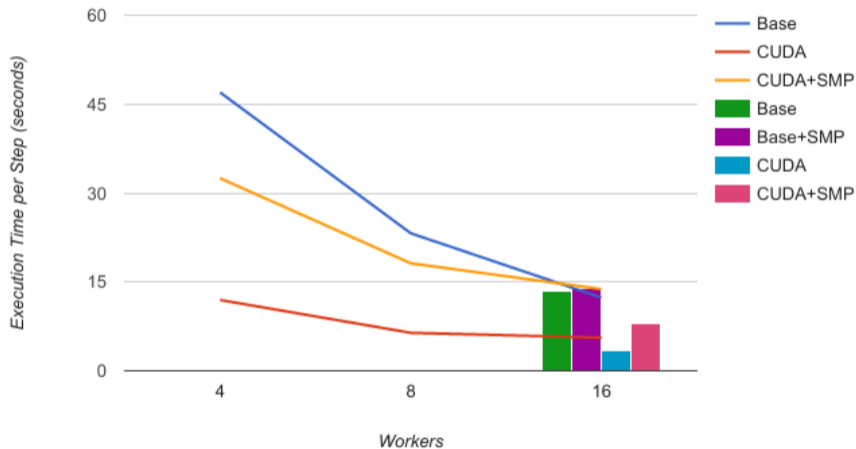


Figure: ChaNGa dwf1 on 4 XK Nodes of BlueWaters

## GPU Applications: ChaNGa GPU Tree Walk

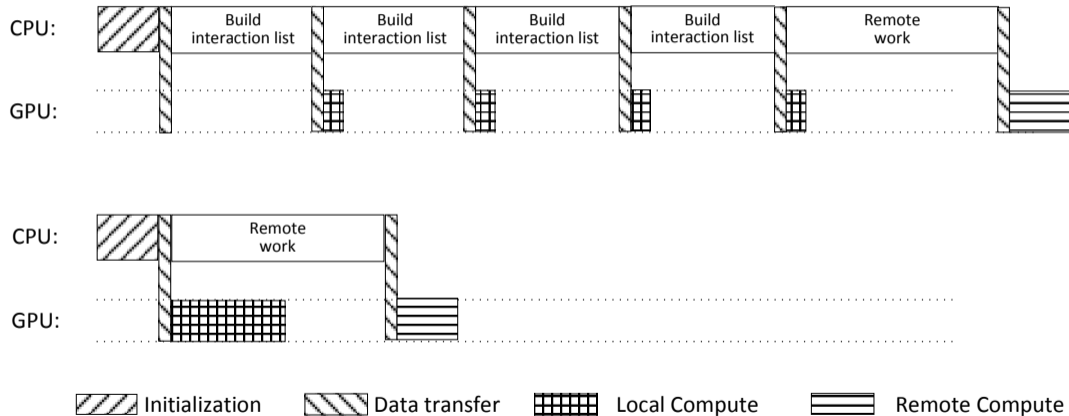


Figure: Strategy Comparison  
Jianqiao Liu, Purdue University

## GPU Applications: ChaNGa GPU Tree Walk

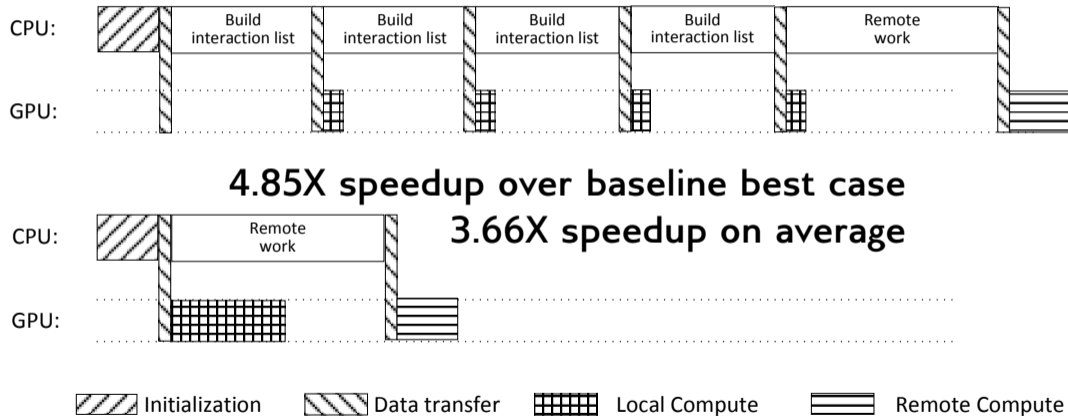
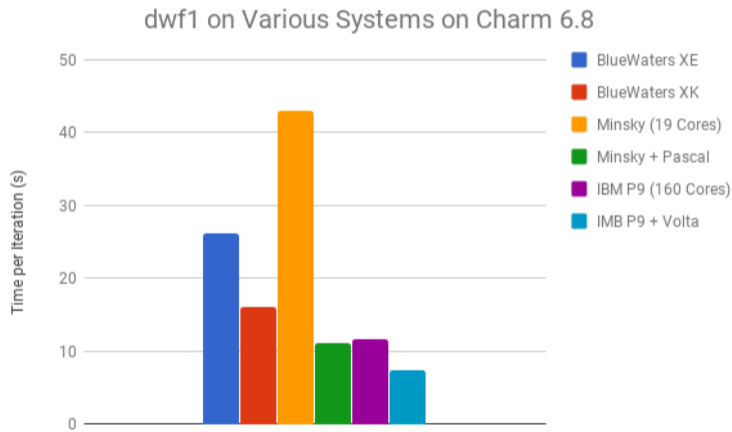


Figure: Strategy Comparison  
Jianqiao Liu, Purdue University

## GPU Applications: ChaNGa on GPU Generations



Mert Hidayetoglu, University of Illinois