# Dynamic Scheduling for Work Agglomeration on Heterogeneous Clusters

Jonathan Lifflander, G. Carl Evans, Anshu Arya, Laxmikant Kale
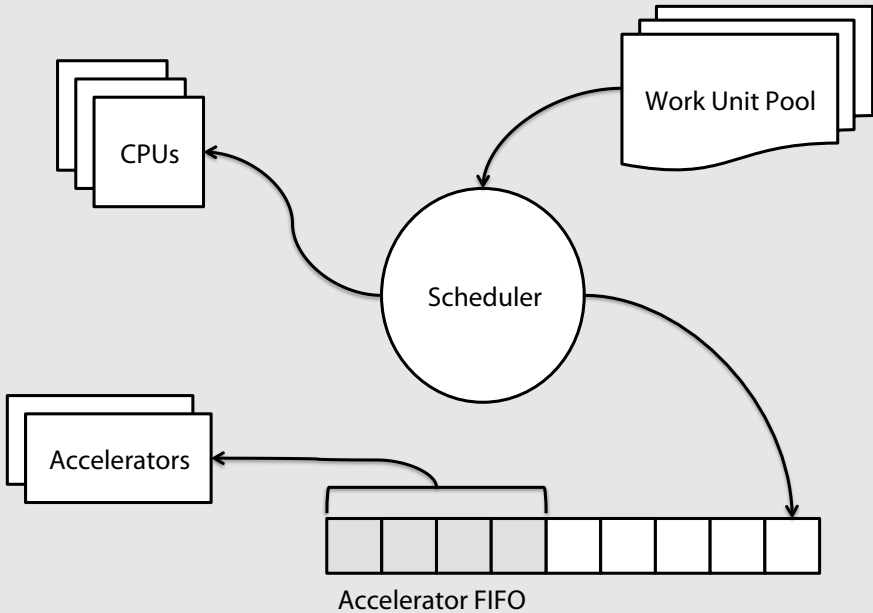
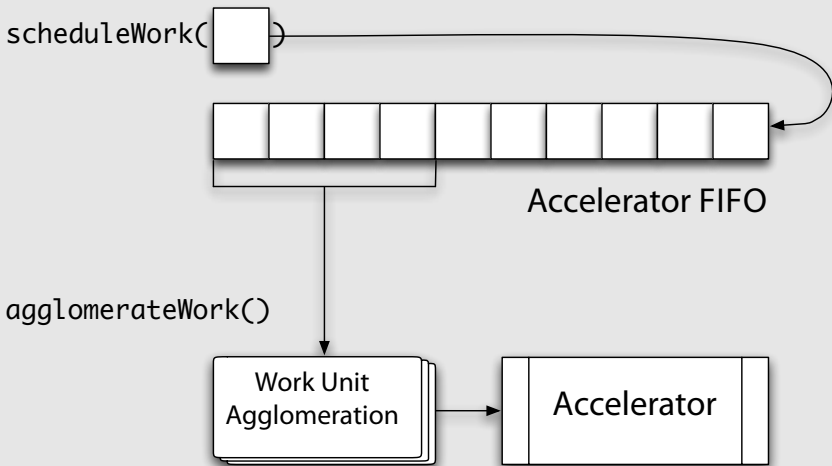University of Illinois Urbana-Champaign

May 7, 2012

- Work is *overdecomposed* in objects
    - Fine-grain task parallelism
    - Ideal for CPU
        - Overlap of communication and computation
    - GPUs rely on massive data-parallelism
        - Fine grains decrease performance
        - Each kernel instantiation has substantial overhead
- To reduce overhead
    - Combine fine-grain work units for the GPU
    - Delay may be insignificant if the work is low priority

# Terminology
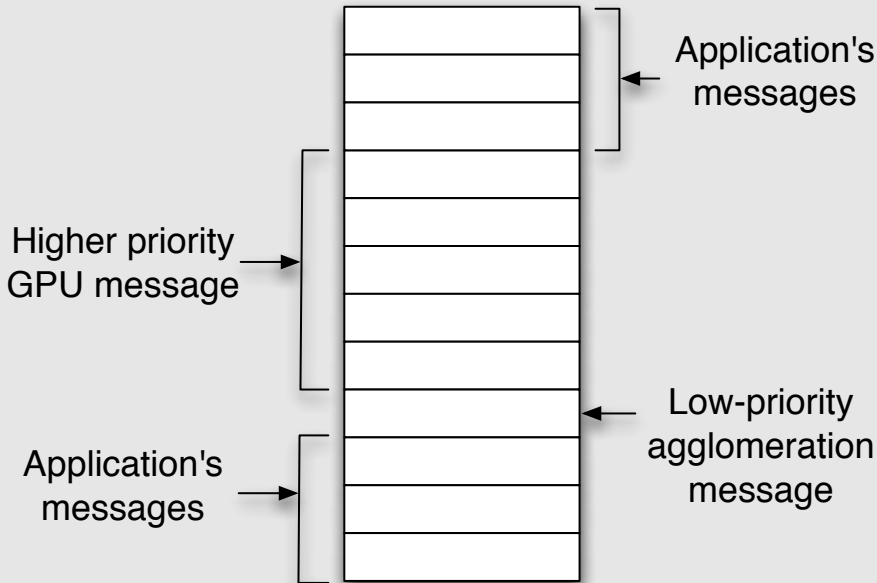
- *Agglomeration*—composition of distinct work units
- *Static agglomeration*—fixed number of work units are agglomerated
- *Dynamic agglomeration*—number of work units agglomerated varies at runtime

CPUs

Work Unit Pool

Scheduler

Accelerators

Accelerator FIFO

scheduleWork( )

Accelerator FIFO

agglomerateWork()

Work Unit
Agglomeration
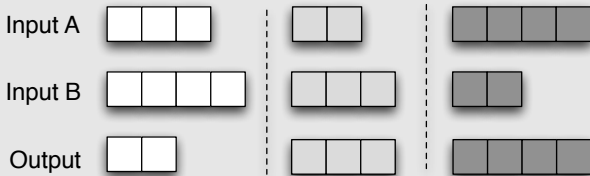
Accelerator

# Programmer/Runtime Division

- ► Programmer
  - ► Writes GPU kernel for agglomeration
  - ► Creates an *offset array*
    - ► Each task's input might be a different size
    - ► Store the offset of each task's beginning and ending index in the contiguous data arrays

- ► System
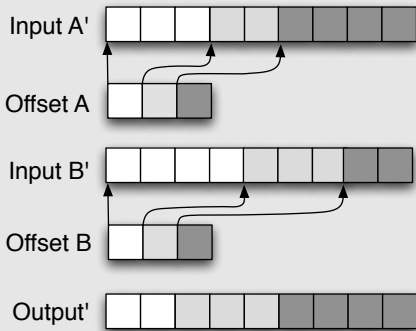  - ► Decide what work to execute and when

Application's messages

Higher priority GPU message

Low-priority agglomeration message

Application's messages

# Dynamic Agglomeration

- Uses the following heuristic
  - If the "accelerator FIFO" reaches a size limit, work is agglomerated
    - Typically set based on memory limitations
  - Else enqueue a low priority message that causes agglomeration
    - When higher-priority work is being generated, it goes into the FIFO
    - When it lets up, work is agglomerated
    - Since low priority work is assumed, not agglomerating aggressively should not impact performance

Non-Agglomerated Data

Input A

Input B

Output

Agglomerated Data
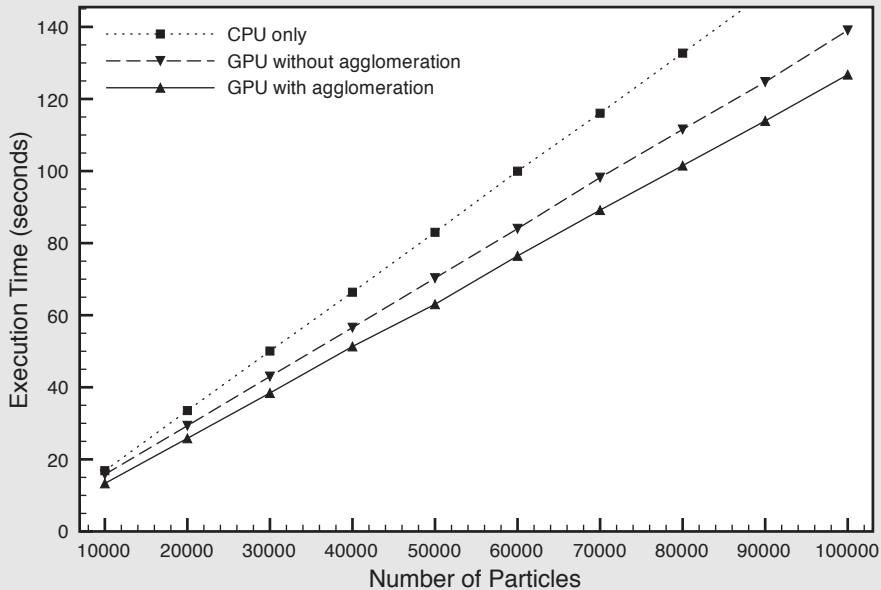
Input A'

Offset A

Input B'
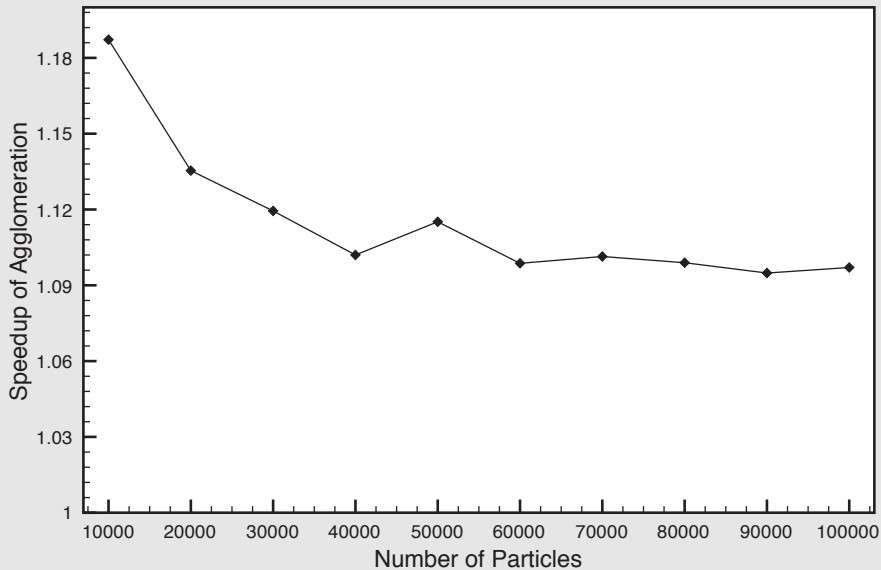
Offset B

Output'

**Case study:** Molecular2D

# Molecular2D
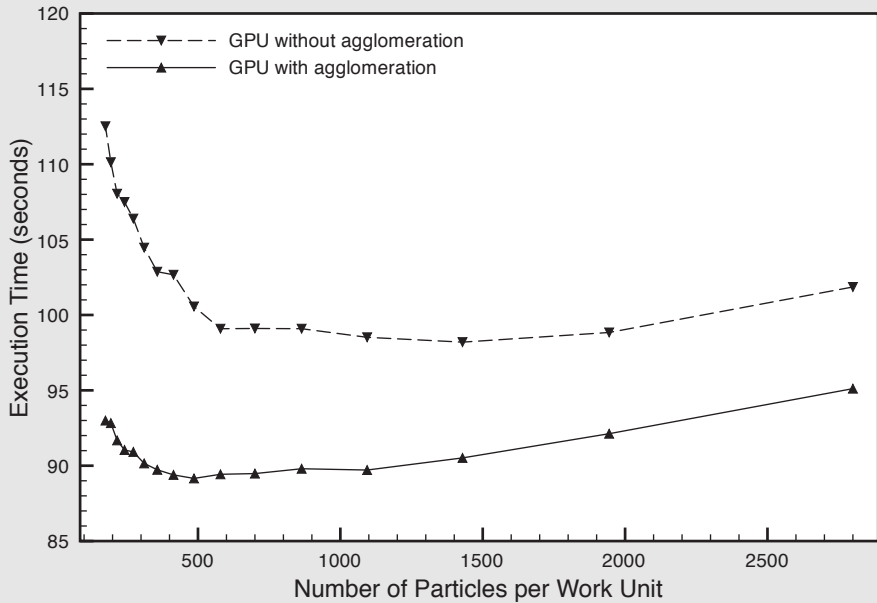
- Cells
    - Execute on CPU
- Interactions
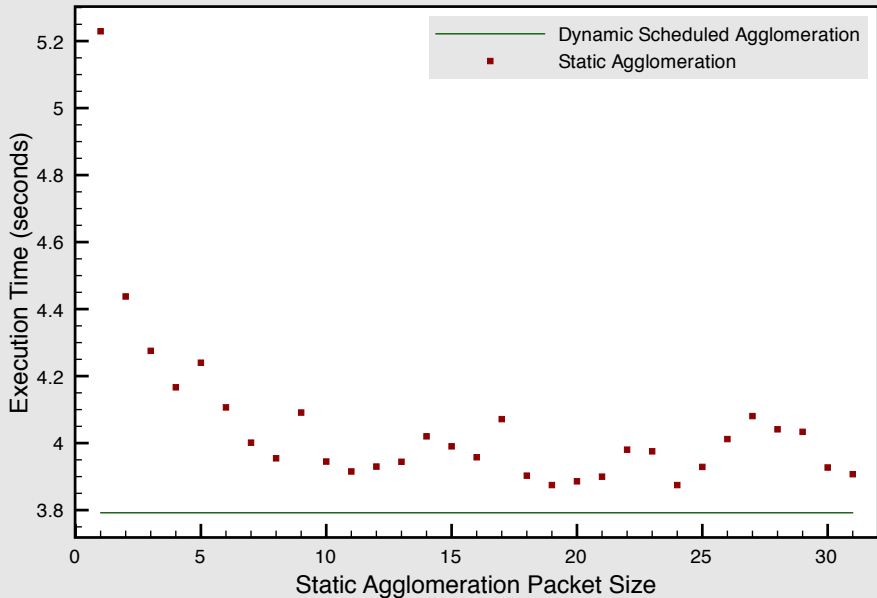    - Execute on GPU

# Molecular 2D Interaction Kernel

```
__global__ void interact(...) {
  int i = blockIdx.x * blockDim.x + threadIdx.x;

  // For loop added for agglomeration
  for(int j = start[i]; j < end[i]; j++)
    // interaction work
}
```
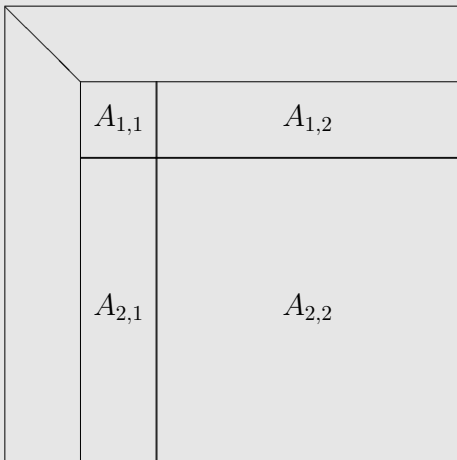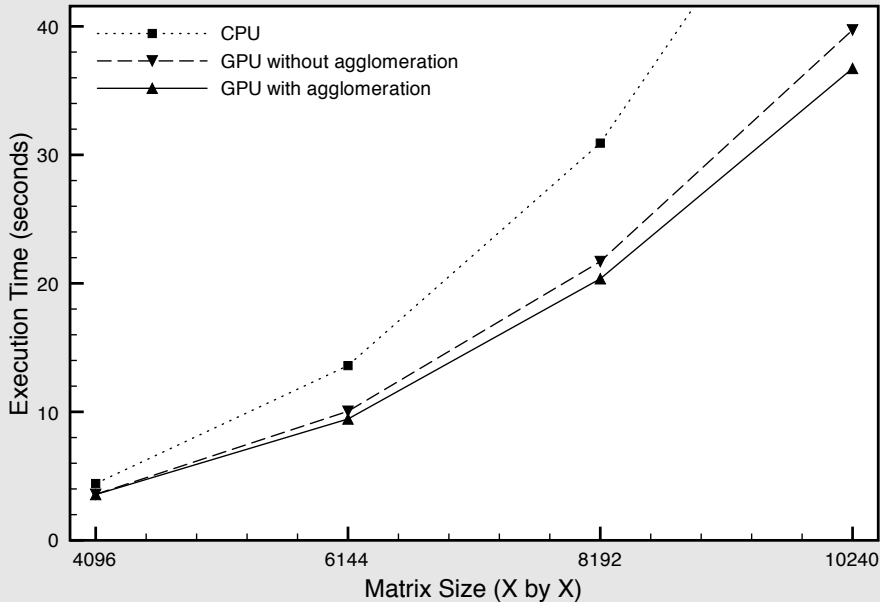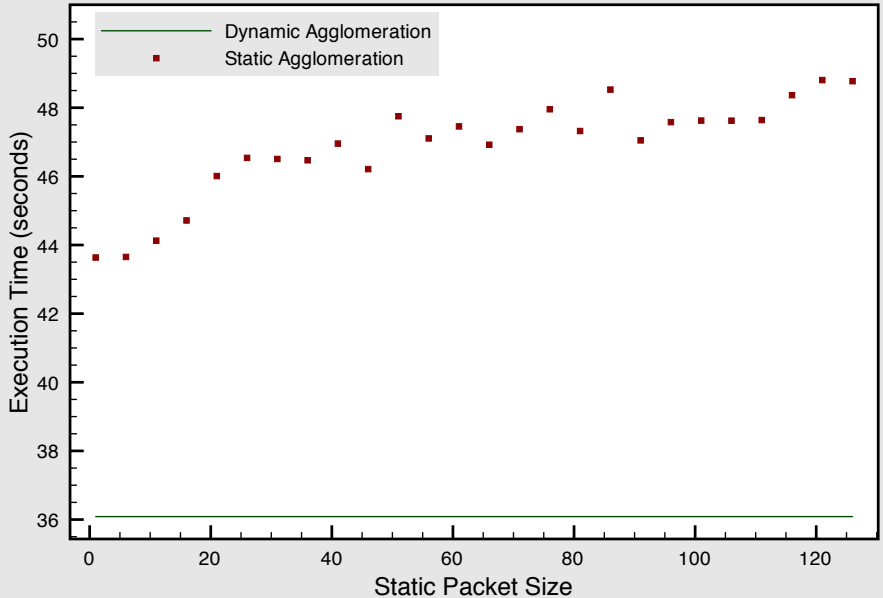
**Case study:** LU Factorization without pivoting

# LU Factorization

- CPU
  - Diagonal
  - Triangular solves
- GPU
  - Matrix-matrix multiples

# Conclusion

- For both benchmarks, agglomerating work increases performance
- Agglomeration does not need to be application-specific
- Statically selecting work units to agglomerate is difficult and may reduce performance
- Runtimes can agglomerate automatically
  - An agglomerating kernel still must written
  - Obtains better performance than static