# Extreme Scale Heterogeneous Computing

**Wen-mei Hwu**

University of Illinois, Urbana-Champaign

# Extreme Scale Heterogeneous Computing Gaining Momentum

- Chinese (Tienhe and Nebula) and Japenese ( Tsubame 2.0) machines rose to top of Top500
- Two of the top 3 of Green500 are Heterogeneous Computing Clusters

# GPU computing is catching on.

| | | | | |
|---|---|---|---|---|
| Financial Analysis | Scientific Simulation | Engineering Simulation | Data Intensive Analytics | Medical Imaging |
| Digital Audio Processing | Digital Video Processing | Computer Vision | Biomedical Informatics | Electronic Design Automation |
| Statistical Modeling | Ray Tracing Rendering | Interactive Physics | Numerical Methods | |

- 280 submissions to GPU Computing Gems
  - 110 articles included in two volumes

# Samples of UIUC Heterogeneous Computing Application Efforts

- NAMD/VMD (Phillips/Stone)
- MILC (Gotlieb/Shi)
- QMC (Kim/Ceperley)
- De Novo Gene Assembly (Ma/Chen/Hwu)
- IMPATIENT (Sutton/Liang/Hwu)
- …

SC 2010

# A Common GPU Usage Pattern

- A desirable approach considered impractical
  - Due to excessive computational requirement
  - But demonstrated to achieve domain benefit
  - Convolution filtering (e.g. bilateral Gaussian filters), De Novo gene assembly, etc.

- Use GPUs to accelerate the most time-consuming aspects of the approach
  - Kernels in CUDA or OpenCL
  - Refactor host code to better support kernels

5
- Rethink the domain problem

# AVAILABLE KERNELS

SC 2010

# Library Kernels

- CUBLAS
  - Basic Linear Algebra
  - CUDA SDK

- CULA, Magma
  - Linear Algebra Solvers
  - www.culatools.com
  - http:icl.cs.utk.edu/magma

- CUSP
  - Sparse data structures and algorithms
  - SpMV, CG, …

- Graph algorithms
  - BFS kernels exist
  - Need graph partitioning kernels

- Unstructured grid algorithms
  - 3D surface mesh generation/ refinement
  - Need 3D volume mesh generation (e.g. CGAL)/ refinement

- Add your favorite library here

# Four Challenges

- Computations with no known scalable parallel algorithms
  - Shortest path, Delaunay triangulation, …
- Data distributions that cause catastrophical load imbalance in parallel algorithms
  - Free-form graphs, MRI spiral scan
- Computations that do not have data reuse
  - Matrix vector multiplication, …
- Algorithm optimizations that are hard and labor intensive
  - Locality and regularization transformations

# Kernel development for GPUs is heavy lifting.

Each kernel is typically a 3-month job but very few developers benefit from advanced compiler technology today.



Little code reuse due to kernel sensitivity to memory access patterns and work distribution.

9

# KERNEL DEVELOPMENT

SC 2010

# Many-core Kernel Development

- Many-core programming is about performance and scalability.
  - Scalability is also key to power efficiency.
  - Performance and scalability for many-cores requires largely the same techniques.
    - To regularize work and data for massively parallel execution.
    - To localize data for conserving memory bandwidth

- There is a gap between what the programmers need and what the tools provide today.

11

# Key to Massive Parallelism - Regularity and Locality

SC 2010

# Eight Optimization Patterns for Algorithms (so far)

| Technique | Contention | Bandwidth | Locality | Efficiency | Load Imbalance | CPU Leveraging |
|---|---|---|---|---|---|---|
| Tiling | | X | X | | | |
| Privatization | X | | X | | | |
| Regularization | | | | X | X | X |
| Compaction | | X | | | | |
| Binning | | X | X | X | | X |
| Data Layout Transformation | X | | X | | | |
| Thread Coarsening | X | X | X | X | | |
| Scatter to Gather Conversion | X | | | | | |

http://courses.engr.illinois.edu/ece598/hk/

GPU Computing Gems, Vol. 1 and 2

# 1: Scatter to Gather Transformation

in

Thread 1    Thread 2    . . .

out

in

Thread 1    Thread 2    . . .

out

# 2. Privatization

Block 0    Block 1    …    Block N

Copy 0    Copy 1    …    Copy N

Final Copy

Parallel Vector Reduction

Block 0    Block 1    …    Block N

🔒    🔒    🔒

Final Copy

🔒 Atomic Updates

15

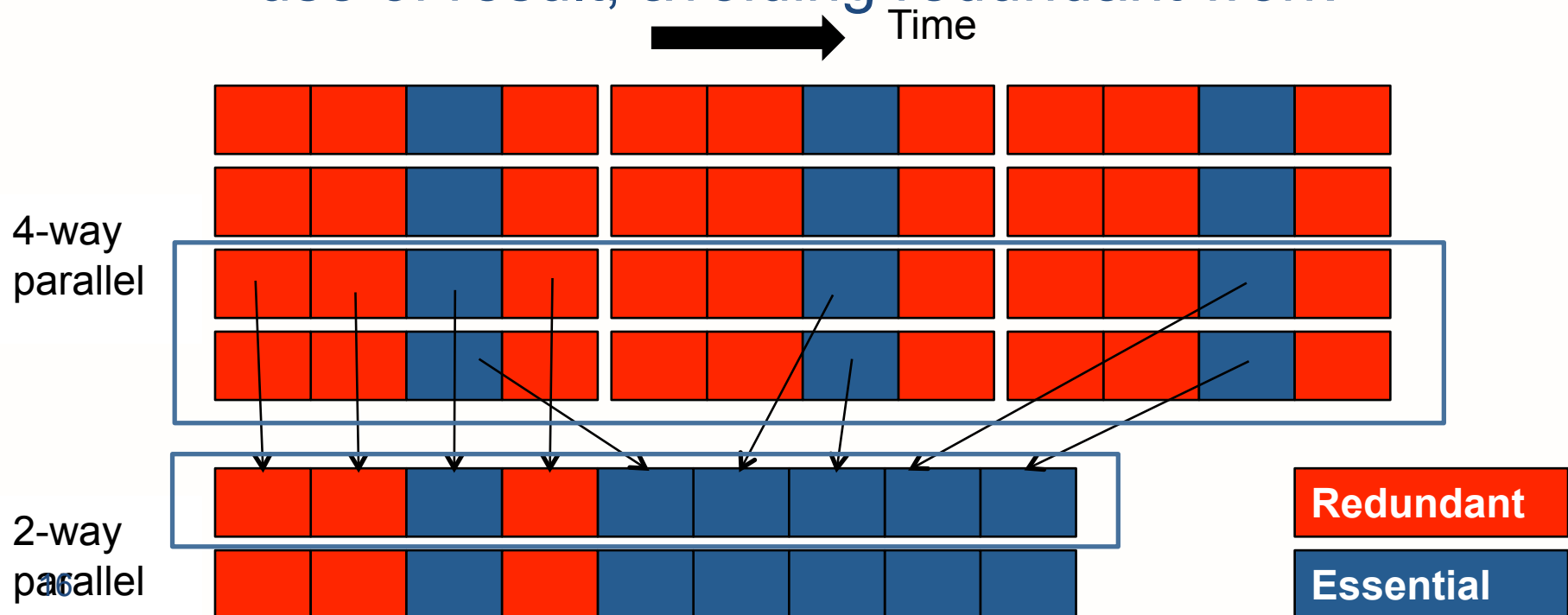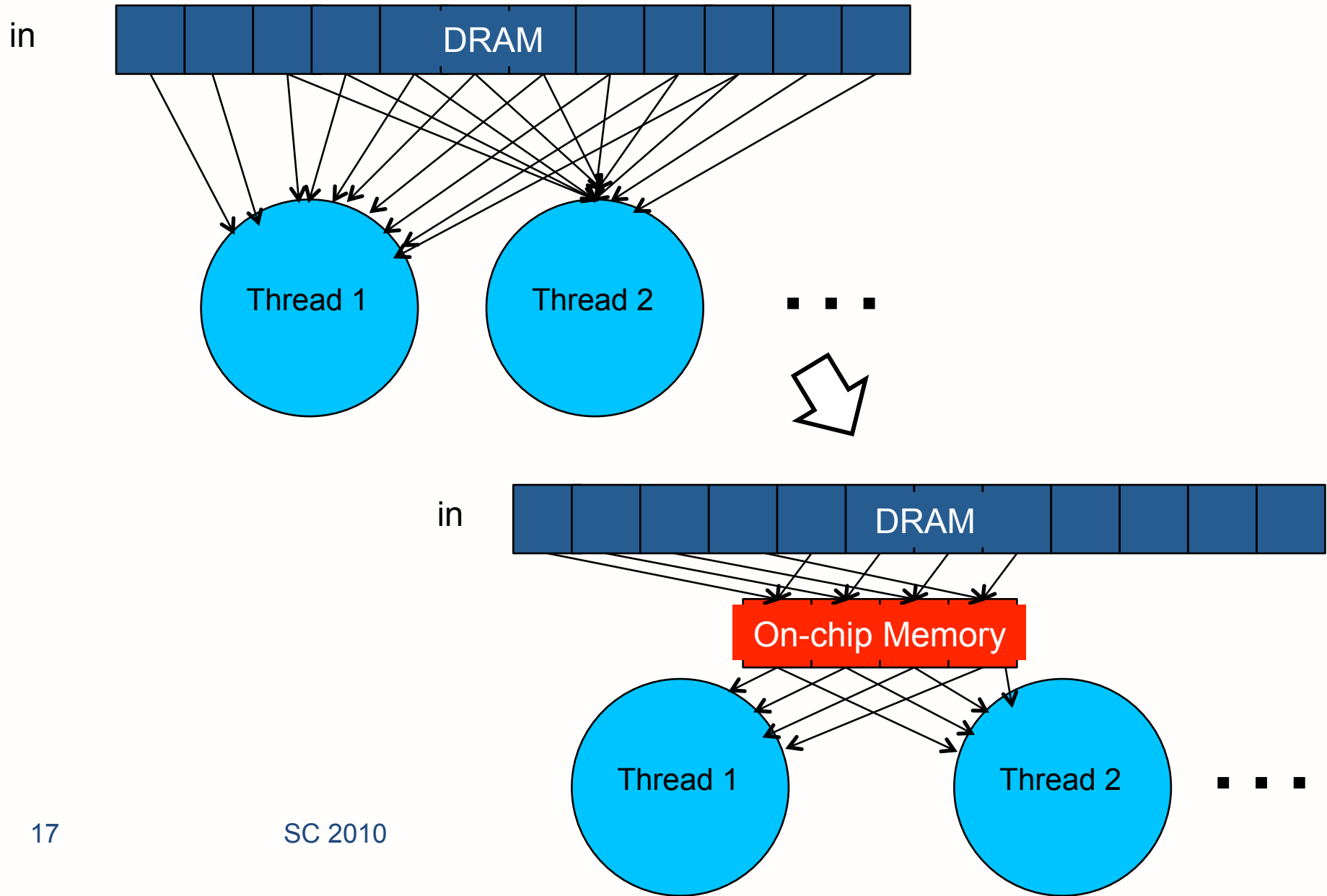# 3. Granularity Coarsening and Register Data Reuse

- Parallel execution often requires redundant and coordination work

  – Merging multiple threads into one allows re-use of result, avoiding redundant work



4-way parallel

2-way parallel

Time

Redundant

Essential

# 4: Data Access Tiling

in | DRAM

Thread 1    Thread 2    . . .

in | DRAM

On-chip Memory

Thread 1    Thread 2    . . .

# 5. Data Layout Transformation



Array of Structure: [z][y][x][e]

Structure of Array: [e][z][y][x]
4X faster than AoS on GTX280

SC 2010 [z][y$_{31:4}$][x$_{31:4}$][e]**[y$_{3:0}$][x$_{3:0}$],** 6.6 X faster than AoS

# 6: Input Data Binning



Bins far beyond the cutoff distance are never scanned



19

# 7. Compaction



Variable sized bins, sort and scan

CPU

On-chip Memory

On-Chip Memory

20

# 8. Regularization

w-queue

On-chip Memory

b-queue

b-queue

g-queue

Work Threads   Dummy Threads

Level i

Propagate

b-queue

Level i+1

b-queue

Level i+2

# Tools go with techniques.

- Tools should facilitate key techniques
  - Programmers should write code "for others to understand instead of for computers to execute" - Dijkstra

- Techniques vary in their potential for automation
  - Scatter-to-gather, granularity coarsening, data access tiling, and memory layout quite amenable
    - Need clear performance guidance
  - Input binning, bin sorting, and hierarchical queues are much harder
    - Need to provide APIs understood by compilers/tools
    - Developer feedback critical to success

# Orion: Reducing Performance Cost of Heterogeneous Parallelism

| CUDA Code | OpenCL Code | Pyon Code | DSL Code |
|-----------|-------------|-----------|----------|

## Orion Performance Portability Framework

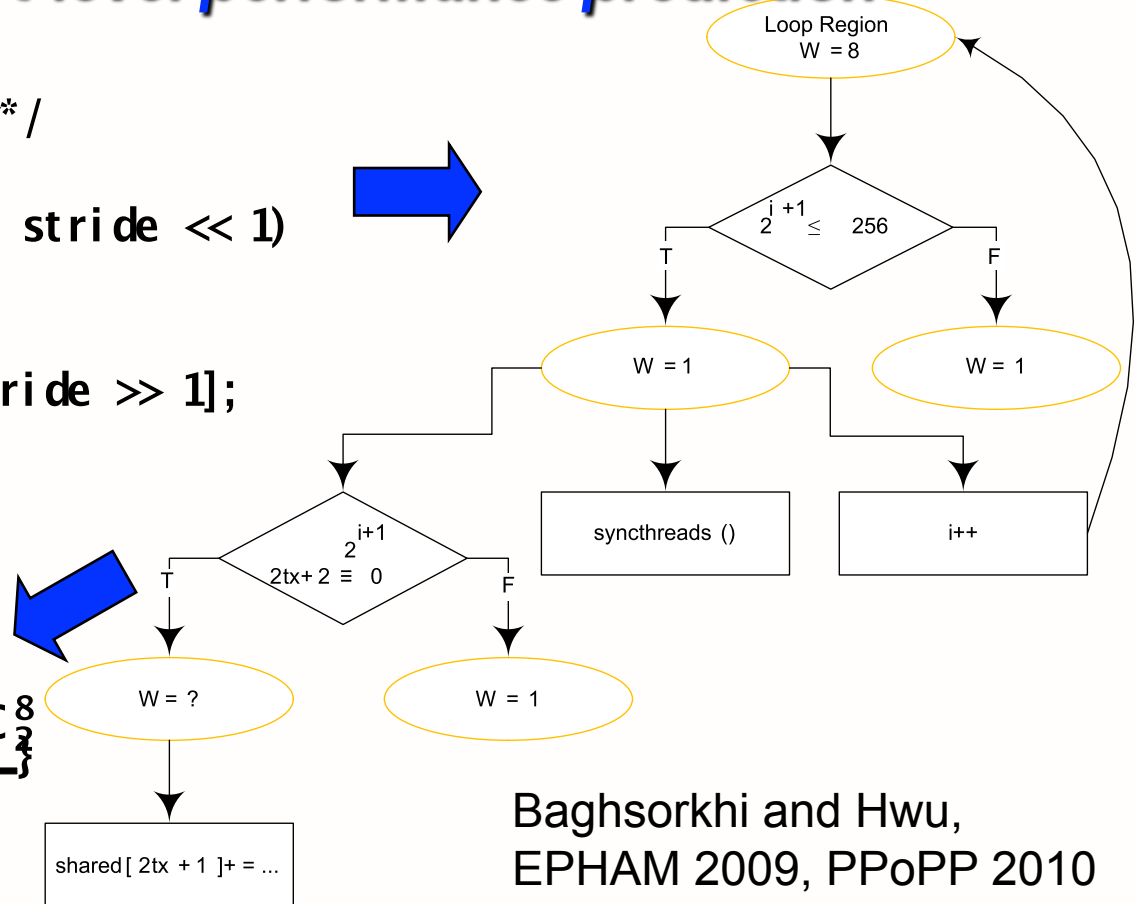| Thread Granularity Coarsening | Tiling and Blocking | Scatter to Gather |
|-------------------------------|---------------------|-------------------|
| Compute and Data Vectorization | Data Layout Transformation | Performance Estimation |

| MuticoreCUDA C Backend | OpenCL Backend | CUDA Backend | FPGACUDA Backend |
|------------------------|----------------|--------------|------------------|

| Host C Compiler | OpenCL Software Stack | NVCC | Auto Pilot Synthesis |
|-----------------|-----------------------|------|----------------------|

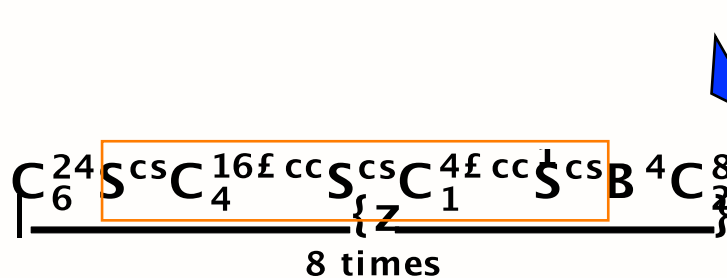| Intel CPUs | AMD CPUs IBM CPUs? | AMD/ATI GPUs | NVIDIA GPUs | Xlinx FPGA |
|------------|--------------------|--------------|-------------|------------|

23

# ADAPT: Example of Advanced Compiler Technique in kernel performance prediction

## HW constraints enable efficient abstract interpretation to emulate expert-level performance prediction
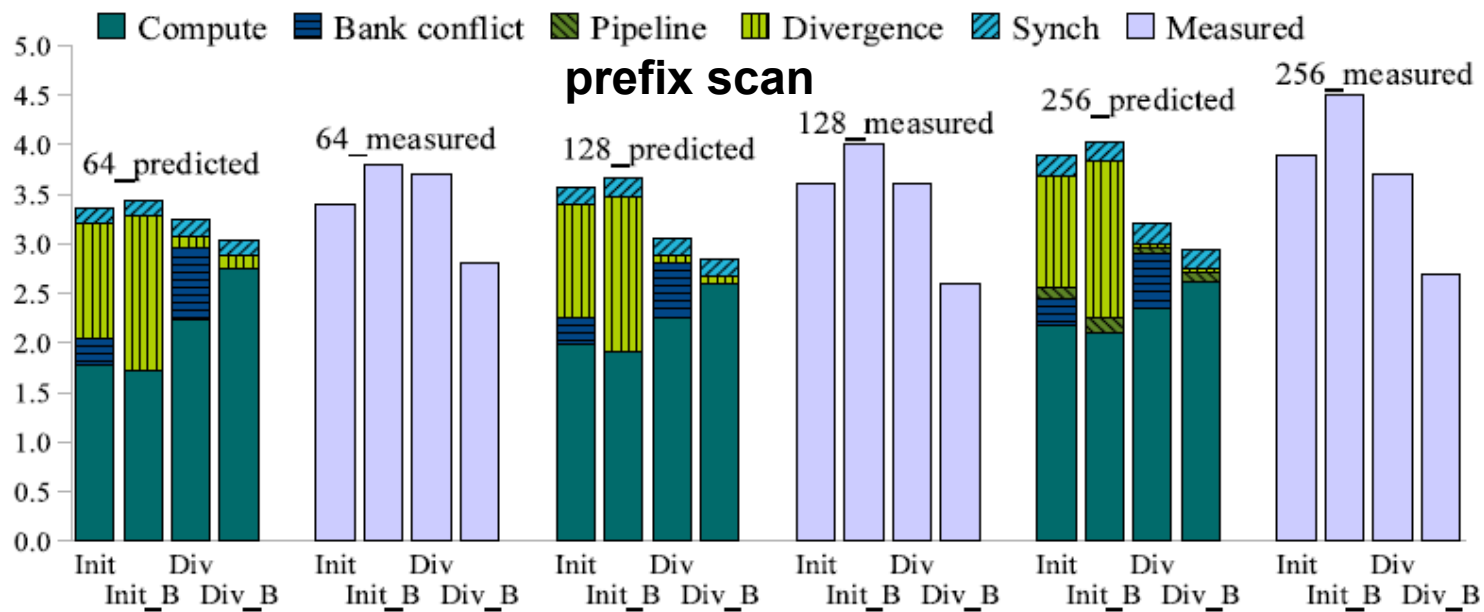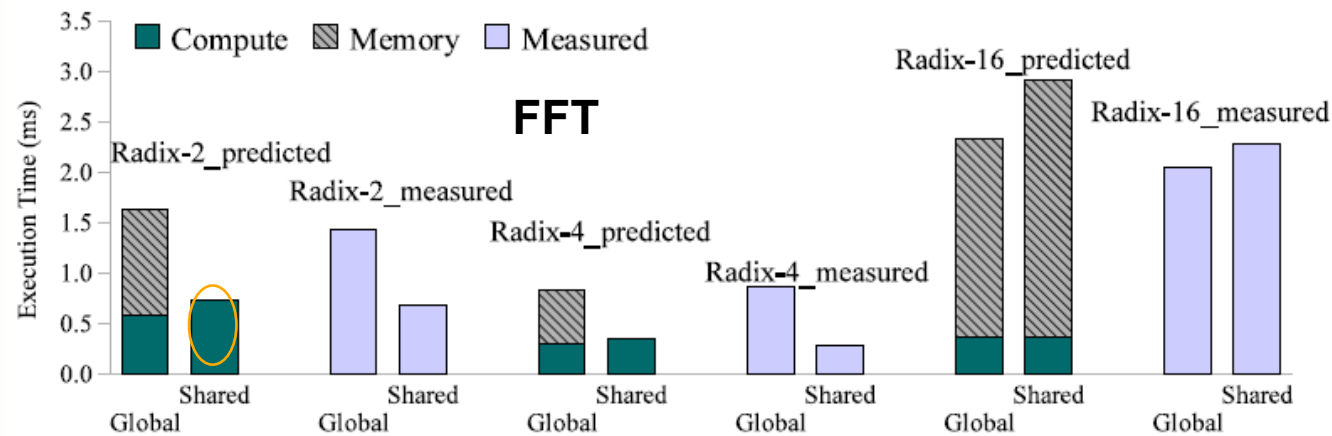
```
n = 2 * tx + 1;
/*Load data into shared memory*/
...
for(stride = 2; stride <= 256; stride << 1)
{
   if( ((n+1) % stride  == 0)
       shared[n] +=shared[n - stride >> 1];
   syncthreads();
}
```
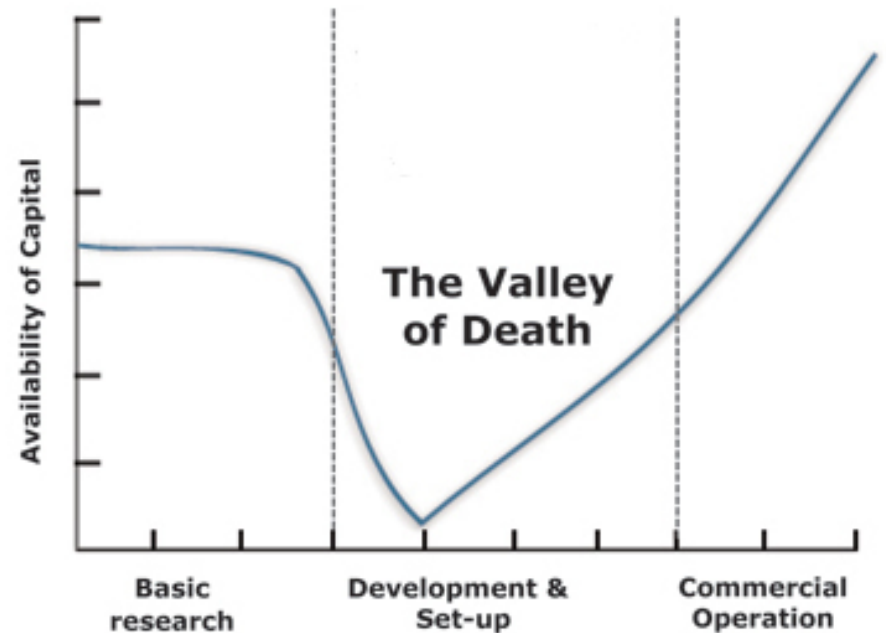


$$C_6^{24} S^{cs} C_4^{16 £ cc} S^{cs} C_1^{4 £ cc} S^{cs} B^4 C_2^8$$

$$\underbrace{\qquad\qquad}_{8 \text{ times}}$$

Loop Region
W = 8

$2^{i+1} \leq 256$

W = 1

W = 1

syncthreads ()

i++

$2^{i+1}$
$2tx+2 \equiv 0$

W = ?

W = 1

shared [ 2tx + 1 ]+ = ...

Baghsorkhi and Hwu,
EPHAM 2009, PPoPP 2010

FFT



prefix scan

# Invitation for Collaboration

- Development and Validation of Scalable Kernel Libraries for Heterogeneous Computing
  - Linear algebra, graph algorithms, PDE solvers, Fourier methods, …
  - New methods/algorithms/implementations
  - Performance portability tools
  - Validation methodology and tools
  - Usable libraries

# Crossing the Valley of Death



We can make it through the valley by collaborating with each other.

# THANK YOU!

SC 2010