# Fast Tridiagonal Solvers on GPU

Yao Zhang

John Owens

UC Davis

Jonathan Cohen

NVIDIA

GPU Technology Conference 2009

# Outline

- Introduction
- Algorithms
  - Design algorithms for GPU architecture
- Performance
  - Bottleneck-based vs. Component-based performance model
- Summary

# What is a tridiagonal system?

$$\begin{pmatrix} b_1 & c_1 & & & & \\ a_2 & b_2 & c_2 & & & \\ & a_3 & b_3 & c_3 & & \\ & & \ddots & \ddots & \ddots & \\ & & & \ddots & \ddots & c_{n-1} \\ & & & & a_n & b_n \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} d_1 \\ d_2 \\ d_3 \\ \vdots \\ \vdots \\ d_n \end{pmatrix}$$

# What is it used for?

- Scientific and engineering computing
  - Alternating direction implicit (ADI) methods
  - Numerical ocean models
  - Semi-coarsening for multi-grid solvers
  - Spectral Poisson Solvers
  - Cubic Spline Approximation
- Video games and computer-animated films
  - Depth of field blurs
  - Fluid simulation

# A Classic Serial Algorithm

- Gaussian elimination in tridiagonal case (Thomas algorithm)

$$\begin{pmatrix} 1 & c'_1 & & & \\ 0 & 1 & c'_2 & & \\ & 0 & 1 & c'_3 & \\ & & 0 & 1 & c'_4 \\ & & & 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{pmatrix} = \begin{pmatrix} d'_1 \\ d'_2 \\ d'_3 \\ d'_4 \\ d'_5 \end{pmatrix}$$
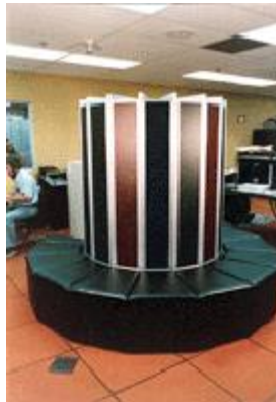
Phase 1: Forward Elimination
Phase 2: Backward Substitution

# Parallel Algorithms

- Coarse-grained algorithms (multi-core CPU)
  - Two-way Gaussian elimination
  - Sub-structuring method

  **A set of equations mapped to one thread**

- Fine-grained algorithms (many-core GPU)
  - Cyclic Reduction (CR)
  - Parallel Cyclic Reduction (PCR)
  - Recursive Doubling (RD)
  - Hybrid CR-PCR algorithm

  **A single equation mapped to one thread**

# A little history

- Parallel tridiagonal solvers since 1960s:
  - Vector machines: Illiac IV, CDC STAR-100, and Cray-1
  - Message passing architectures: Intel iPSC, and Cray T3E
  - And GPU as well!

# Two Applications on GPU





Depth of field blur, Michael Kass et al.

OpenGL and Shader language
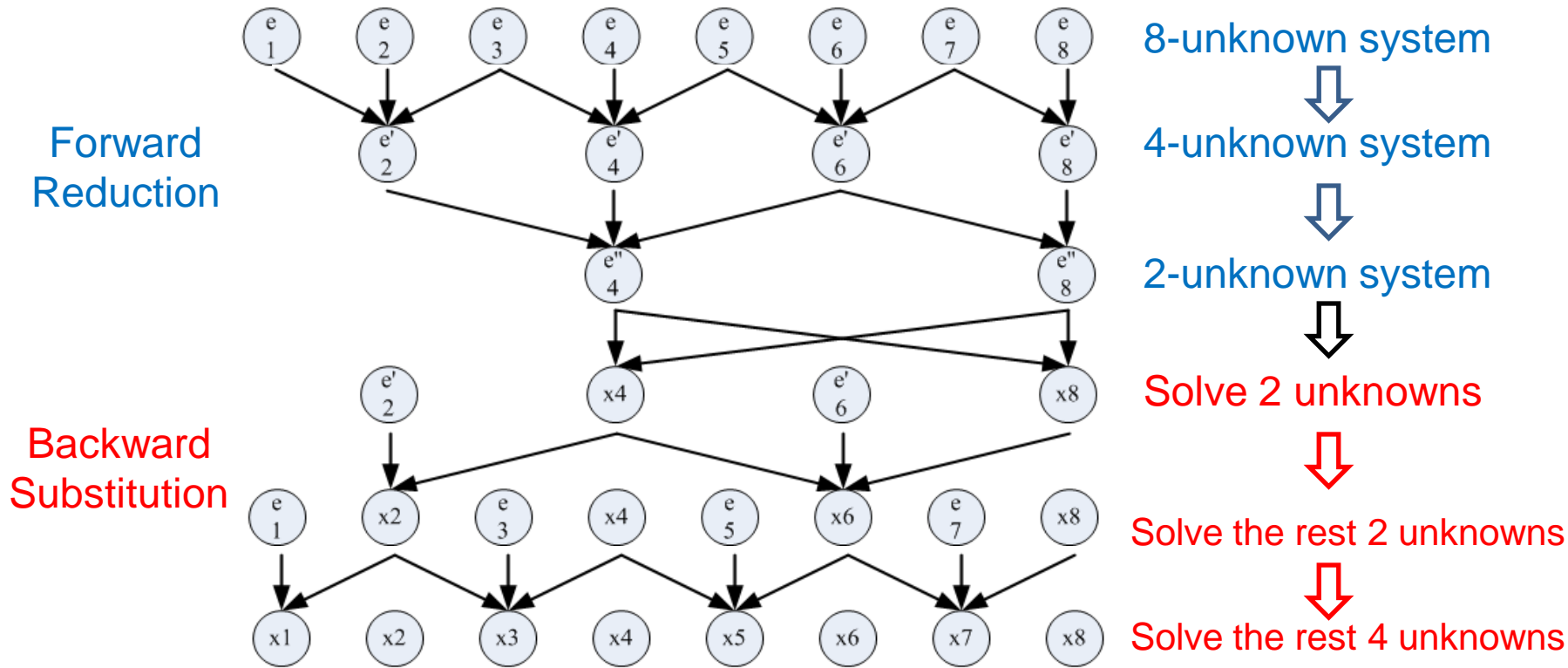
Cyclic reduction

2006

Shallow water simulation
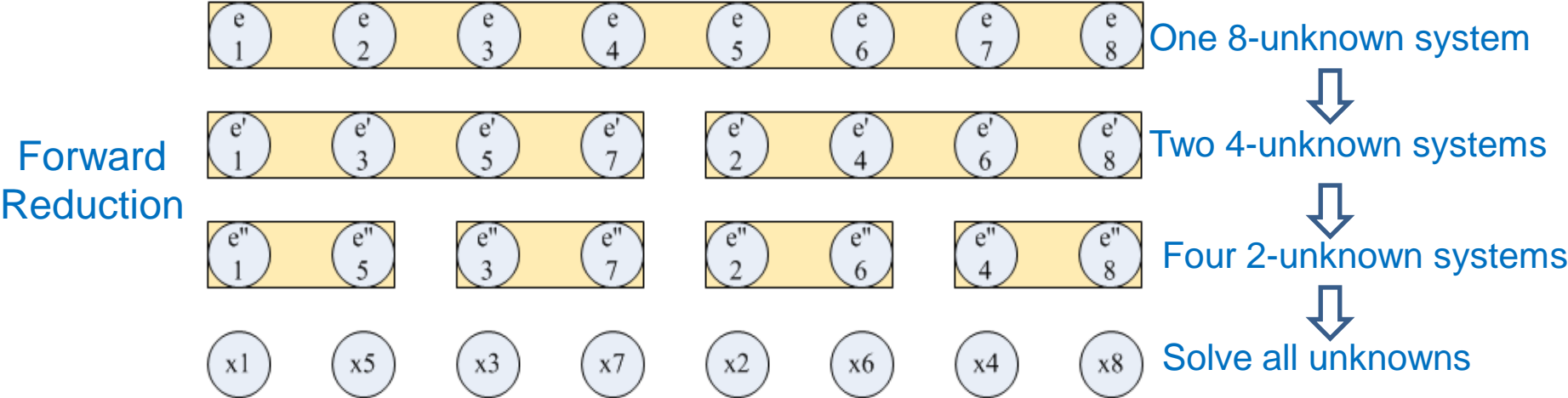
CUDA

Cyclic reduction

2007

# Cyclic Reduction (CR)

# Parallel Cyclic Reduction (PCR)

**4** threads working

Forward Reduction

One 8-unknown system

Two 4-unknown systems
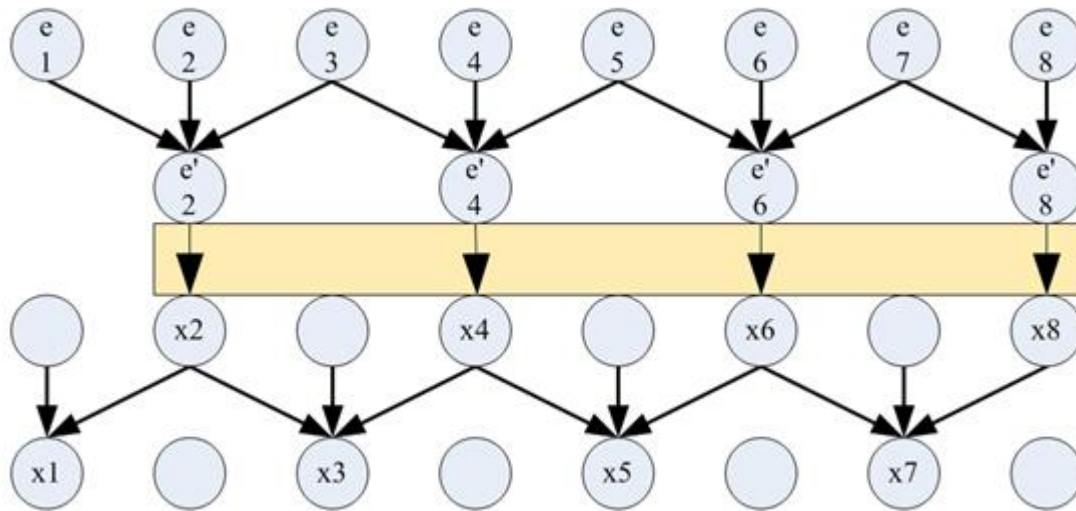
Four 2-unknown systems

Solve all unknowns

$\log_2 (8) = 3$ steps

# Hybrid Algorithm (1)

- CR
  - Every step we reduce the system size by half (Good)
  - Some processing cores stay **idle** if the system size is smaller than the number of cores (Bad)
  - Needs more steps to finish (Bad)
- PCR
  - Fewer steps required (Good)
  - Same amount of work for all steps (Bad)
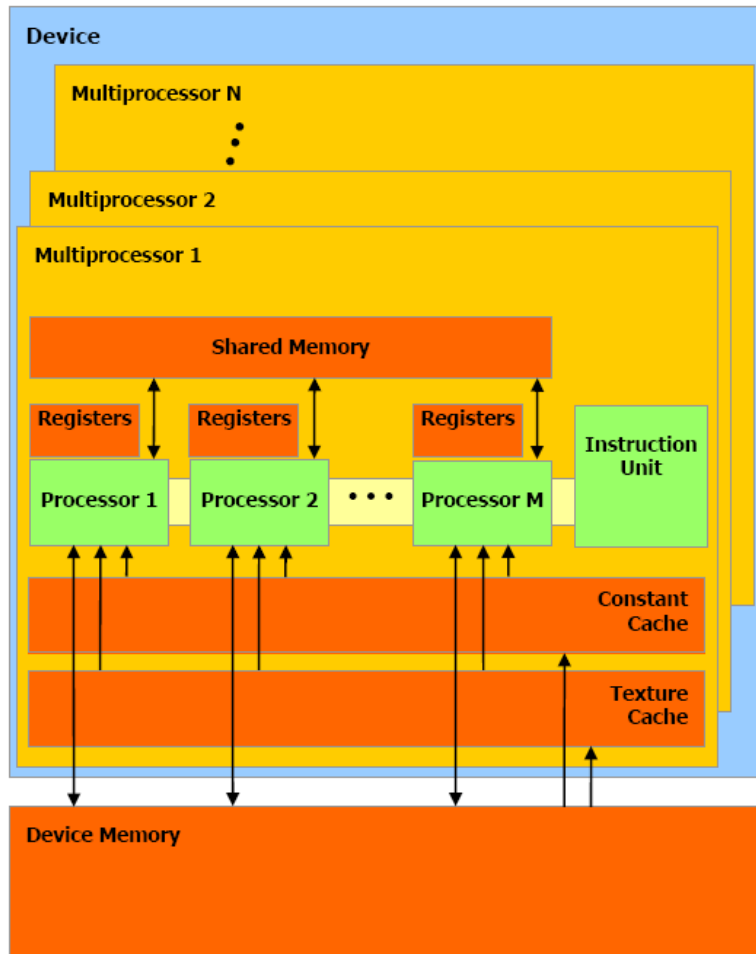
# Hybrid Algorithm (2)



Switch to PCR

Switch back to CR

System size reduced at the beginning
No idle processors
Fewer algorithmic steps

Even more beneficial because of:
bank conflicts
control overhead

# GPU Implementation (1)



- Linear systems mapped to multiprocessors (blocks)

- Equations mapped to processors (threads)
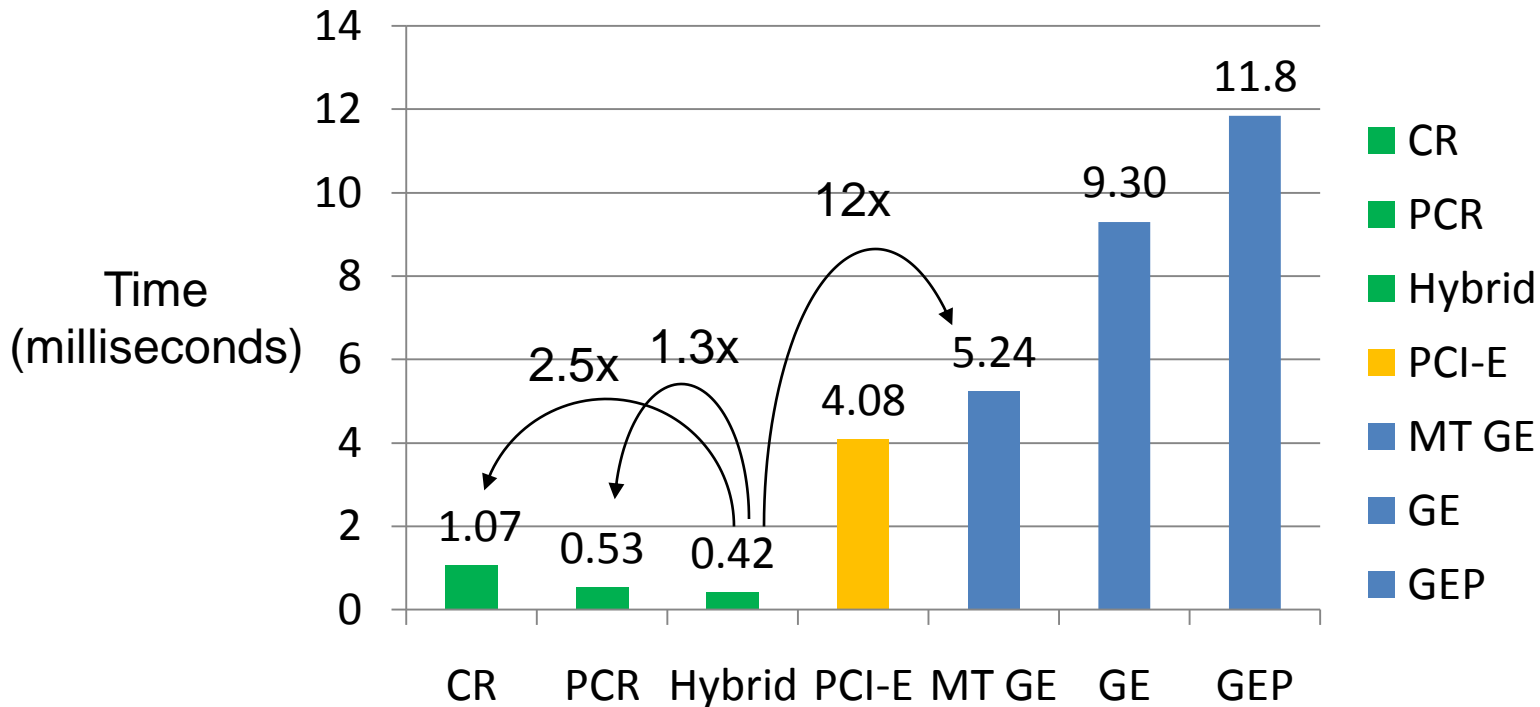
# GPU Implementation (2)

- Storage need: 5 arrays = 3 diagonals + 1 solution vector + 1 right hand side

- All data resides in shared memory if it fits

- Use contiguously ordered threads to avoid unnecessary divergent branches

- In-place data storage
  - Efficient, but introduce bank conflicts to CR

# Performance Results – Test Platform

- 2.5 GHz Intel Core 2 Q9300 quad-core CPU
- GTX 280 graphics card with 1 GB video memory
- CUDA 2.0
- CentOS 5 Linux operating system

# Performance Results



**Solve 512 systems of 512 unknowns**

Time (milliseconds)

Legend: CR, PCR, Hybrid, PCI-E, MT GE, GE, GEP

Values: CR 1.07, PCR 0.53, Hybrid 0.42, PCI-E 4.08, MT GE 5.24, GE 9.30, GEP 11.8

Annotations: 2.5x, 1.3x, 12x

PCI-E: CPU-GPU data transfer
MT GE: multi-threaded CPU Gaussian Elimination
GEP: CPU Gaussian Elimination with pivoting (from LAPACK)

# Performance Analysis

- Factors that determine performance
  - Global/shared memory accesses
  - Bank conflicts
  - Computational complexity
  - Overhead for synchronization and loop control
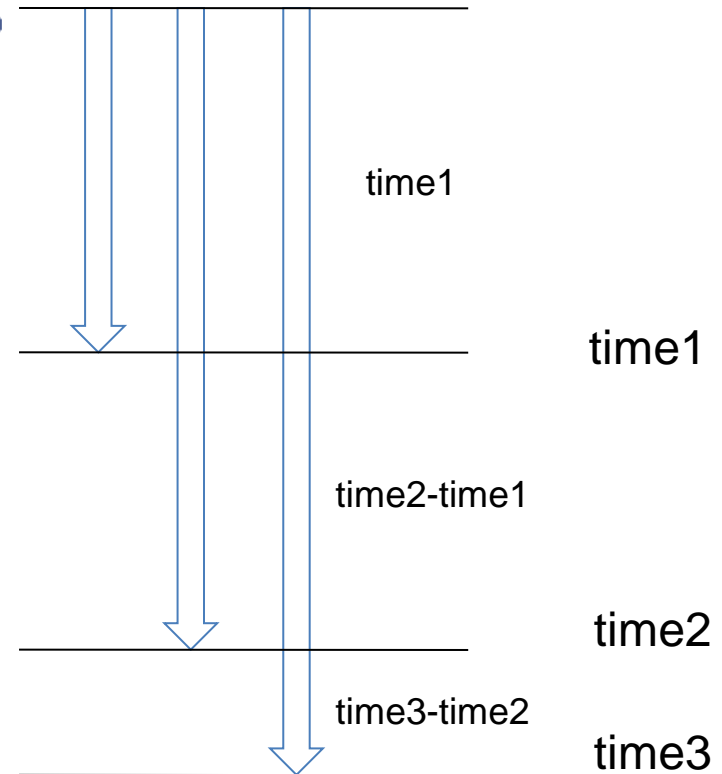
# Bottleneck vs. Pie slice



Performance = **min**(factor1, factor2, …)     Performance = **sum**(factor1, factor2, …)
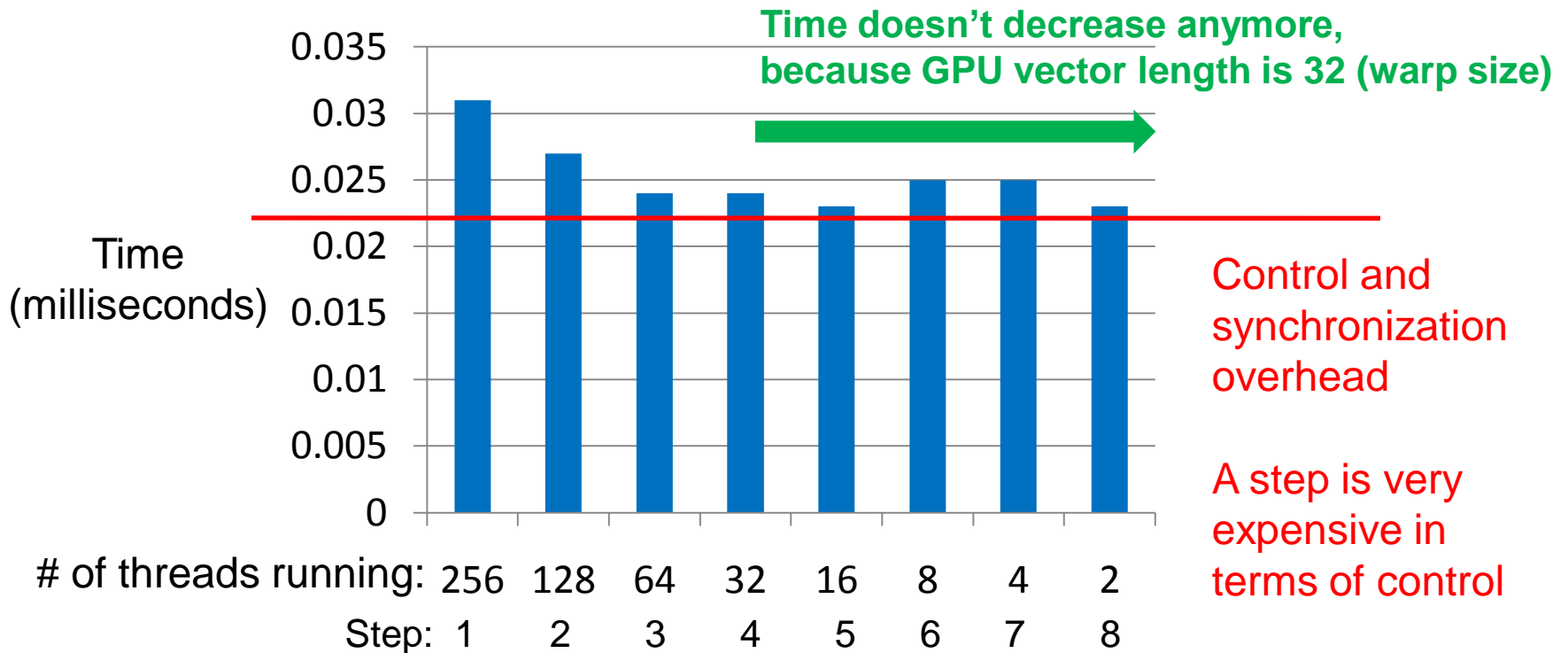
# Performance Measure

A manual differential method

```
46 __global__ void scan_naive(float *g_odata, float *g_idata, int n)
47 {
48     // Dynamically allocated shared memory for scan kernels
49     extern __shared__ float temp[];
50
51     int thid = threadIdx.x;
52
53     int pout = 0;
54     int pin = 1;
55
56     // Cache the computational window in shared memory
57     temp[pout*n + thid] = (thid > 0) ? g_idata[thid-1] : 0;
58
59     for (int offset = 1; offset < n; offset *= 2)
60     {
61         pout = 1 - pout;
62         pin  = 1 - pout;
63         __syncthreads();
64
65         temp[pout*n+thid] = temp[pin*n+thid];
66
67         if (thid >= offset)
68             temp[pout*n+thid] += temp[pin*n+thid - offset];
69     }
70
71     __syncthreads();
72
73     g_odata[thid] = temp[pout*n+thid];
74 }
```
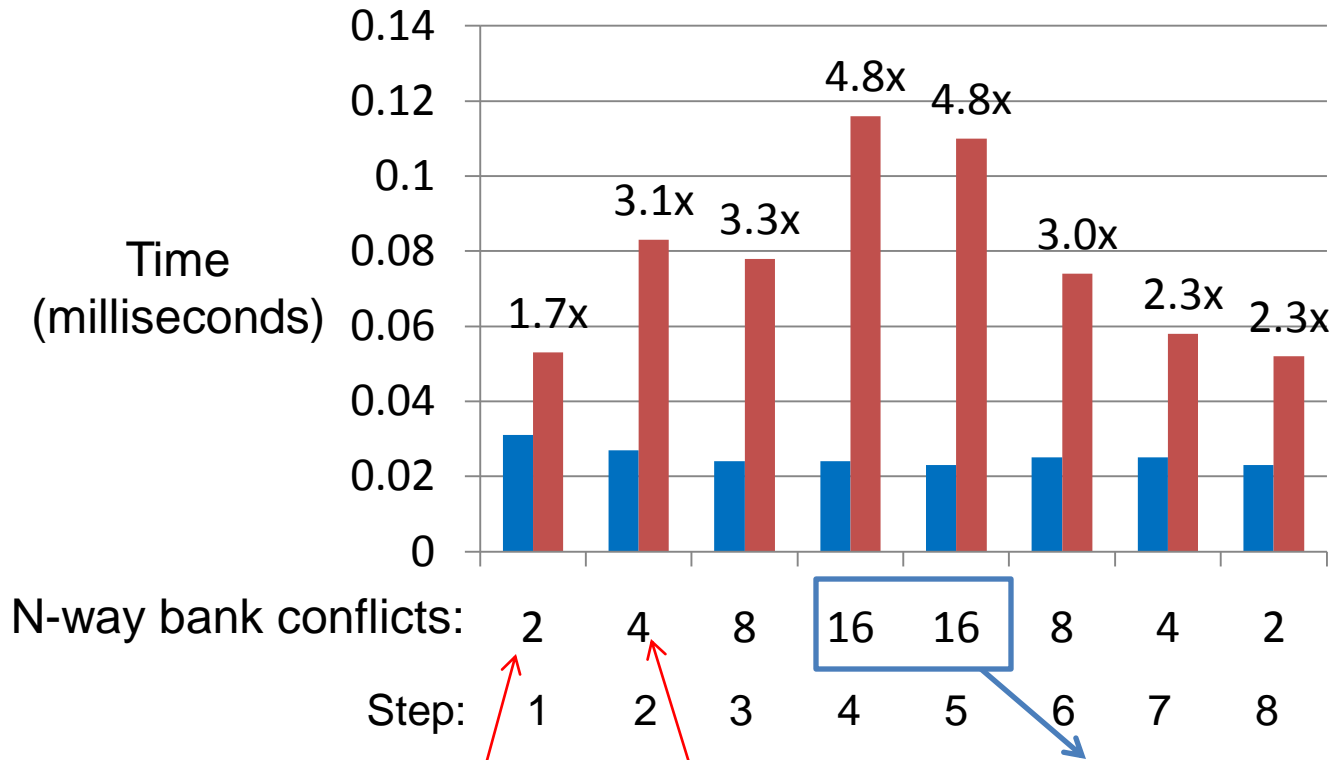
time1

time1

time2-time1

time2

time3-time2

time3

# Control Overhead

**CR: Forward Reduction**



**Time doesn't decrease anymore, because GPU vector length is 32 (warp size)**

Time (milliseconds)

Control and synchronization overhead

A step is very expensive in terms of control

# of threads running: 256  128  64  32  16  8  4  2

Step:  1  2  3  4  5  6  7  8

Enforce a stride of one to avoid bank conflicts

# Bank Conflicts



CR: Forward Reduction

# CR vs. PCR (1)



Solve 512 systems of 512 unknowns
(Time Breakdown)

PCR  0.53 ms
CR   1.07 ms

Forward Reduction    Backward Substitution

Global memory    Solve 2-unknown systems

# CR vs. PCR (2)



**CR**

0.1
9%
49 GB/s

0.27
26%
16 GFLOPS

0.69
65%
33 GB/s

Global  Shared  Computation

**PCR**

0.11
20%
47 GB/s

0.27
50%
102 GFLOPS

0.16
30%
883 GB/s

Global  Shared  Computation

Pros: O(n)
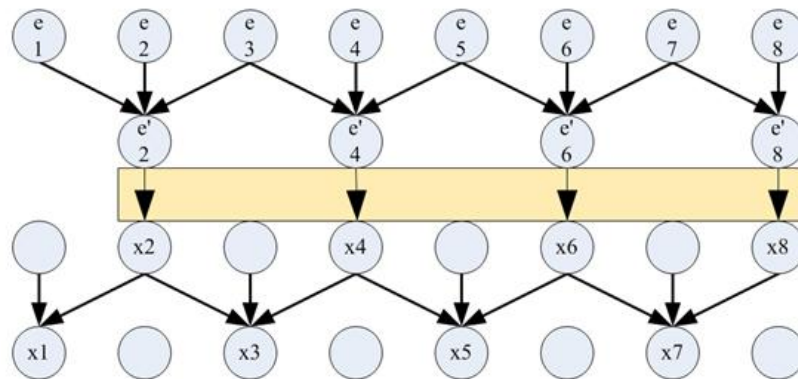Cons: more steps (control overhead), bank conflicts

Pros: fewer steps, no bank conflicts
Cons: O(nlogn)

# Pitfalls

- The higher computation rate and sustained bandwidth, the better
  - They may have different algorithm complexity
- The lower algorithm complexity, the better
  - What if there is considerable amount of control overhead, or bank conflicts, or low hardware utilization
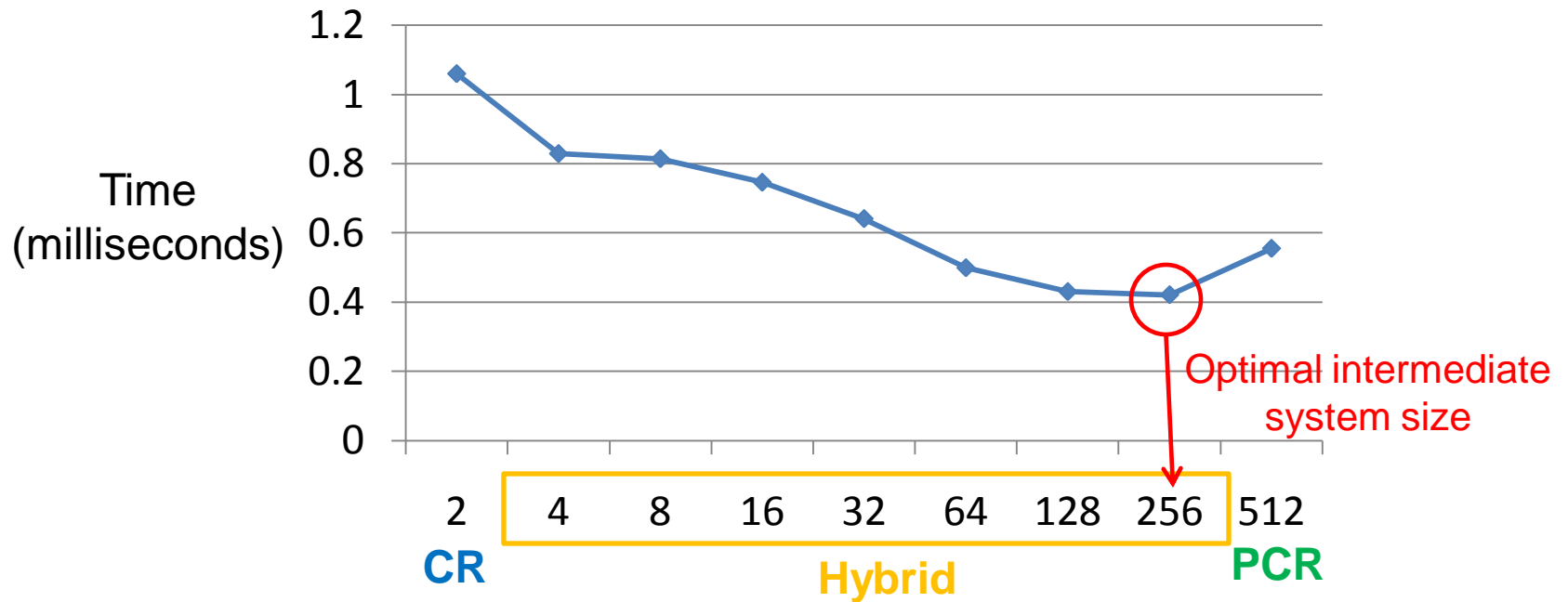
# PCR vs Hybrid

- Make tradeoffs between the computation, memory access, and control
  - The earlier you switch from CR to PCR
    - The fewer bank conflicts, the fewer algorithmic steps
    - But more work

# Hybrid Solver – Sweet Point



Optimal performance of hybrid solver
Solving 512 systems of 512 unknowns

# Known issues and future research

- PCI-E data transfer

- Double precision

- Pivoting

- Block tridiagonal systems

- Handle large systems that cannot fit into shared memory

- Automatic performance profiling

# Summary

- We studied the performance of 3 parallel solvers on GPU

- We learned two major lessons
  - Component-based rather than bottleneck-based
    - » Performance is more complicated than either compute-bound or memory-bound
  - We can make tradeoffs between these components, and we need to make the right tradeoff

# Questions?

Thanks