# A Memory-Efficient Algorithm for Large-Scale Symmetric Tridiagonal Eigenvalue Problem on Multi-GPU Systems

Hyunsu Cho and Peter A. Yoon

Trinity College, Hartford, CT, USA

# Symmetric Eigenvalue Problem

$$A\mathbf{x} = \lambda\mathbf{x}$$

where $A$ is symmetric

Many interesting applications require eigenvectors

# Divide and Conquer

Yields **full spectrum** of eigenvalues and eigenvectors

Is numerically stable

Gives rise to **independent subproblems**

Often faster than $O(n^3)$ due to deflation

# Divide and Conquer

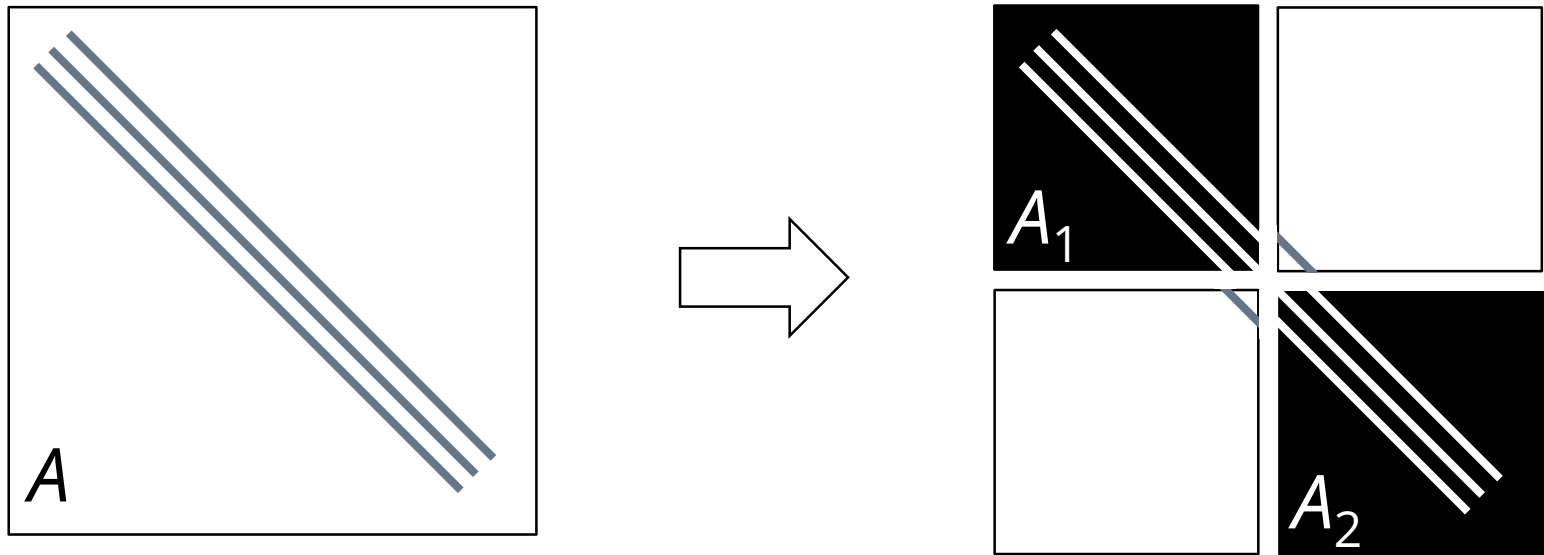Apply **orthogonal similarity transformation** to reduce $A$ to tridiagonal form

$$Q^T A Q = A'$$

where

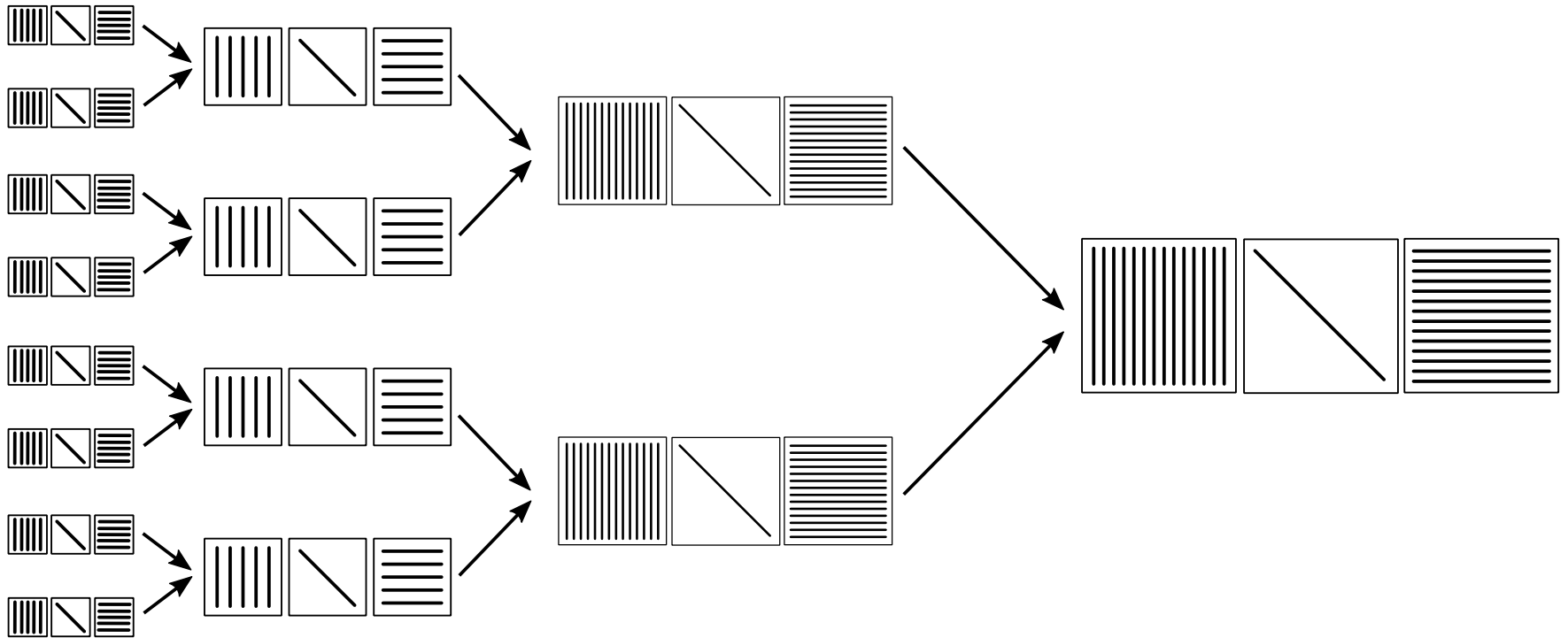$A'$ is symmetric tridiagonal
and $Q$ is orthogonal

Existing work on single-node, multi-GPU:
  MAGMA (UTK)

# Divide and Conquer



- Solve subproblems
- Merge solutions
- Repair

# Divide and Conquer

# Merging solutions

Suppose

$$A = \left[\begin{array}{c|c} A_1 & \\ \hline & A_2 \end{array}\right] + \left[\begin{array}{c|c} & \\ b_m & b_m \\ \hline b_m & b_m \\ & \end{array}\right]$$

where

$$A_1 = Q_1 D_1 Q_1^T \qquad \text{(subproblem \#1)}$$
$$A_2 = Q_2 D_2 Q_2^T \qquad \text{(subproblem \#2)}$$

# Merging solutions

Then

$$A = QDQ^T + \begin{bmatrix} & & b_m & b_m & \\ \hline & & b_m & b_m & \end{bmatrix}$$

where

$$Q = \begin{bmatrix} Q_1 & \\ \hline & Q_2 \end{bmatrix} \quad \text{and} \quad D = \begin{bmatrix} D_1 & \\ \hline & D_2 \end{bmatrix}$$

# Merging solutions

$$H = b_m \left[ \frac{\mathbf{e}_m}{\mathbf{e}_1} \right] \left[ \frac{\mathbf{e}_m}{\mathbf{e}_1} \right]^T$$

Then

$$A = QDQ^T + \begin{bmatrix} \begin{array}{c|c} b_m & b_m \\ \hline b_m & b_m \end{array} \end{bmatrix}$$

where

$$Q = \begin{bmatrix} \begin{array}{c|c} Q_1 & \\ \hline & Q_2 \end{array} \end{bmatrix} \quad \text{and} \quad D = \begin{bmatrix} \begin{array}{c|c} D_1 & \\ \hline & D_2 \end{array} \end{bmatrix}$$

# Rank-one update

$$H = b_m \left[ \frac{\mathbf{e}_m}{\mathbf{e}_1} \right] \left[ \frac{\mathbf{e}_m}{\mathbf{e}_1} \right]^T$$

$$A = QDQ^T + H$$
$$= Q(D + b_m \mathbf{z}\mathbf{z}^T)Q^T$$

where

$$\mathbf{z} = Q^T \left[ \frac{\mathbf{e}_m}{\mathbf{e}_1} \right] = \left[ \begin{array}{c} \text{last column of } Q_1^T \\ \text{first column of } Q_2^T \end{array} \right]$$

# Rank-one update

$$H = b_m \begin{bmatrix} \mathbf{e}_m \\ \hline \mathbf{e}_1 \end{bmatrix} \begin{bmatrix} \mathbf{e}_m \\ \hline \mathbf{e}_1 \end{bmatrix}^T$$

$$A = QDQ^T + H$$

$$= Q(D + b_m \mathbf{z}\mathbf{z}^T)Q^T$$

where

$$\mathbf{z} = Q^T \begin{bmatrix} \mathbf{e}_m \\ \hline \mathbf{e}_1 \end{bmatrix} = \begin{bmatrix} \text{last column of } Q_1^T \\ \text{first column of } Q_2^T \end{bmatrix}$$

Need eigen-decomposition of inner system

# Decompose $D + b_m \mathbf{z} \mathbf{z}^T$

1. Sort entries in $D$; permute z  likewise
2. Filter some entries in $D$ and z via deflation (next slide)

# Decompose $D + b_m \mathbf{z}\mathbf{z}^T$

1. Sort entries in $D$; permute z likewise
2. Filter some entries in $D$ and z via deflation (next slide)
3. Compute all roots of the **secular equation** [1]

$$1 + b_m \sum_{i=1}^{n} \frac{d_i^2}{z_i - \lambda} = 0,$$

   giving the $m$ eigenvalues.
4. Compute corresponding eigenvectors stably [2]

[1] Li 1994
[2] Gu & Eisenstat 1994

# Decompose $D + b_m \mathbf{z}\mathbf{z}^T$

1. Sort entries in $D$; permute z likewise

2. Filter some entries in $D$ and z via deflation (next slide)

3. Compute all roots of the **secular equation** [1]

$$1 + b_m \sum_{i=1}^{n} \frac{d_i^2}{z_i - \lambda} = 0,$$

   giving the $m$ eigenvalues.

4. Compute corresponding eigenvectors stably [2]

5. Multiply each eigenvector by $Q$
   Recall: $A = Q(D + b_m \mathbf{z}\mathbf{z}^T)Q^T$

[1] Li 1994
[2] Gu & Eisenstat 1994

# Deflation

Recall:
$$D = \left[\begin{array}{c|c} D_1 & \\ \hline & D_2 \end{array}\right]$$

Entries of $D$ are eigenvalues of two subproblems

If two entries are nearly identical, we throw one away

**Fewer columns** when multiplying eigenvectors by $Q$

Same thing for small entries in z

**Reduce work complexity to $O(n^{2.3})$**

# GPU computing

**General-purpose** computation on GPUs

**Bulk parallelism** w/ many small threads

Cost effective; widely available

# Mapping work to GPU

1. Sort entries in $D$; permute z likewise

2. Filter some entries in D and z via deflation

3. Compute all roots of the secular equation, giving the $m$ eigenvalues.

4. Compute corresponding eigenvectors stably

5. Multiply each eigenvector by $Q$
   → Done in bulk via DGEMM

Parallel but
not as work-intense

# GPU memory

High-bandwidth dedicated memory
Separate from main memory
Limited in size

| Main memory | GPU memory |
| CPU | GPU |

PCI-E

# Memory requirement

Eigenvectors are dense
$\rightarrow O(n^2)$ storage

Intermediate workspace: eigenvectors of inner system

| Matrix dimension | Memory required |
|---:|---:|
| 8192 | 1.5 GB |
| 16384 | 5.8 GB |
| 32768 | 23.4 GB |
| 36000 | 28.2 GB |
| 50000 | 54.4 GB |

# Our contribution

**Overcome limitation in GPU memory**
while retaining adequate performance

# Strategies

## 1. Use multiple GPUs

# Strategies

2. Keep most of workspace in main memory (**out-of-core** approach)

# Strategies

3. **Shape work** to fit GPU workspaces

# Block matrix multiplication

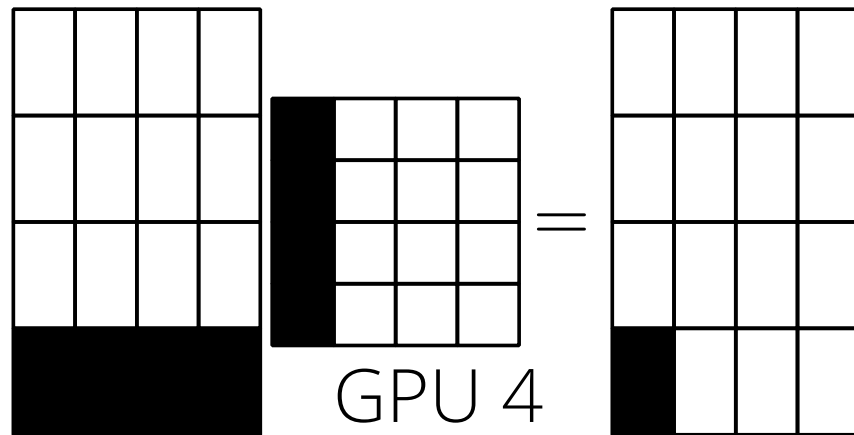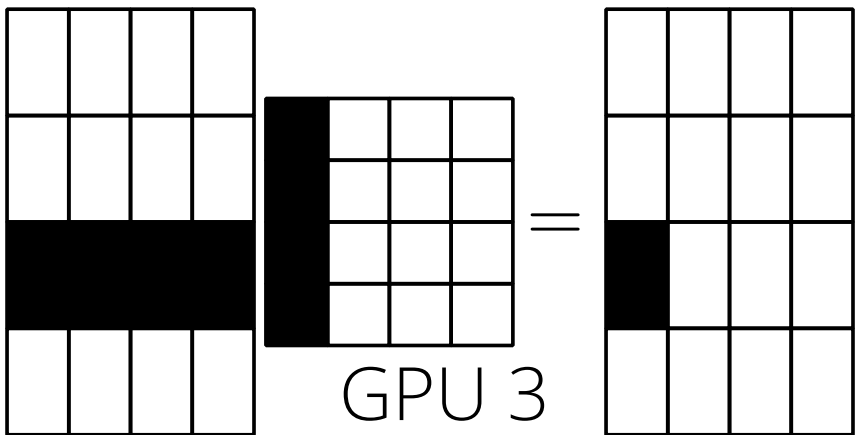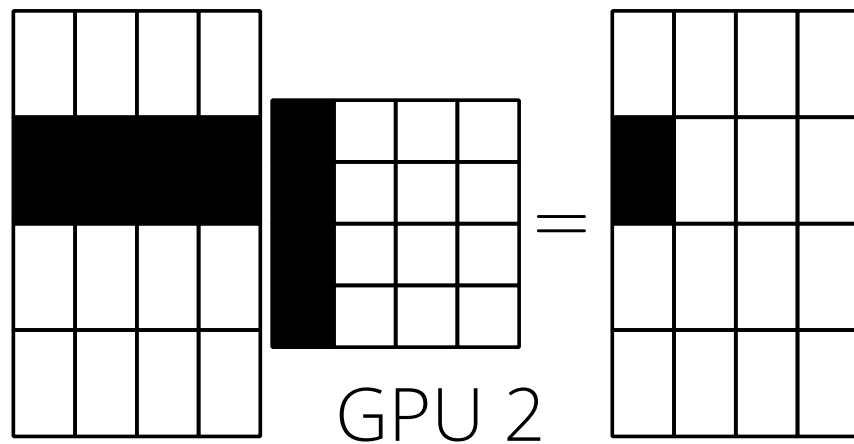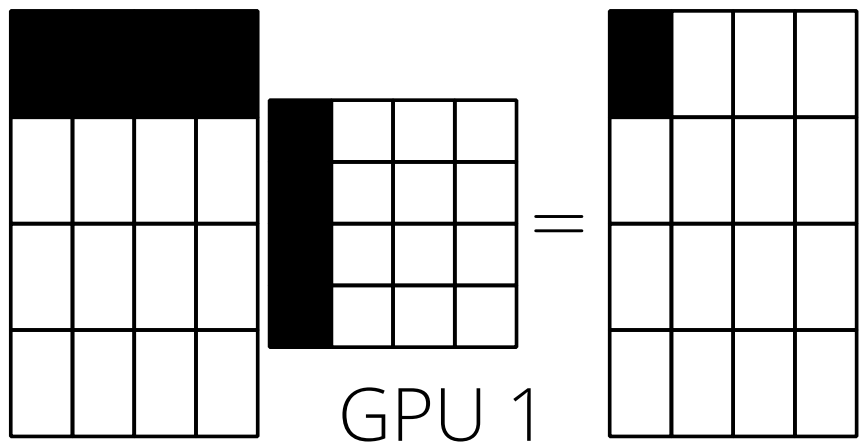Use a fine partition to fit submatrices into GPU memory



$$Q$$

Eigenvectors of
$$D + b_m \mathbf{z}\mathbf{z}^T$$

Eigenvectors of $A$

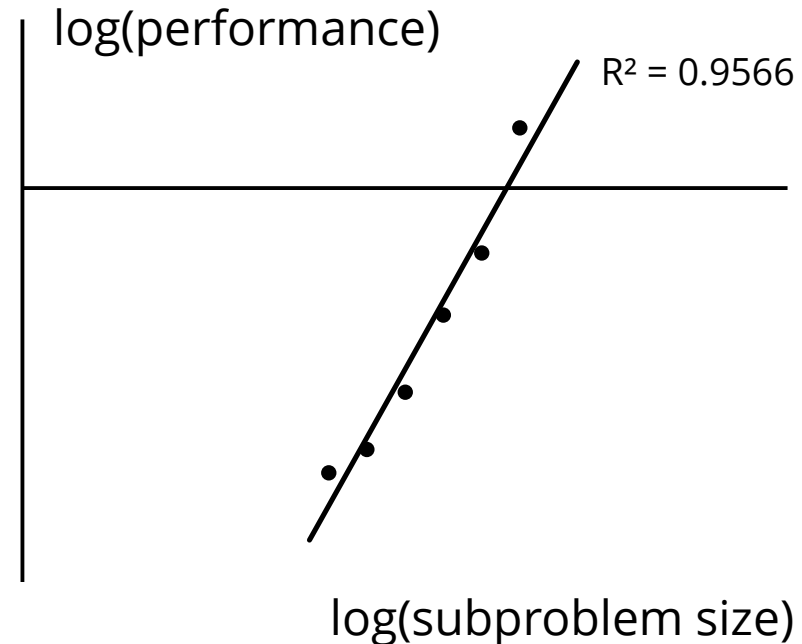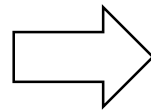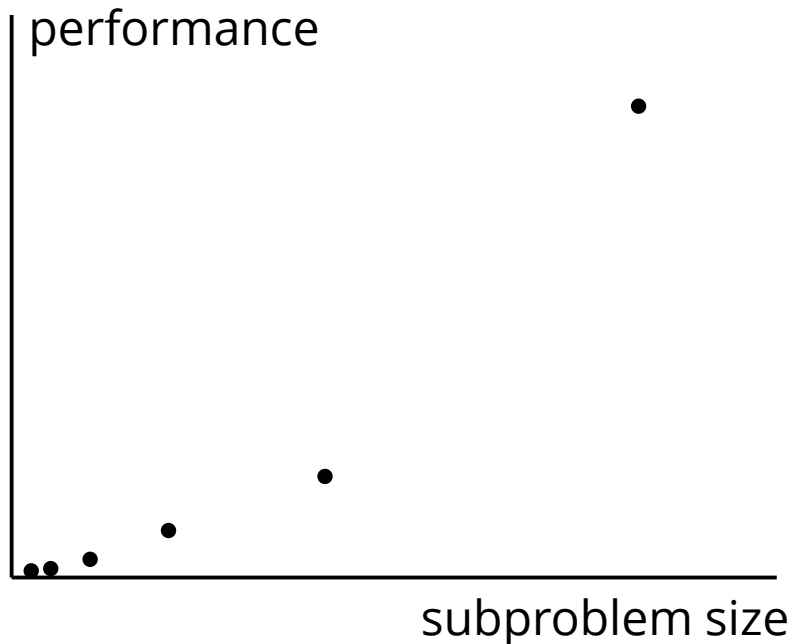GPU 1

GPU 2

GPU 3

GPU 4

# Hybrid computation
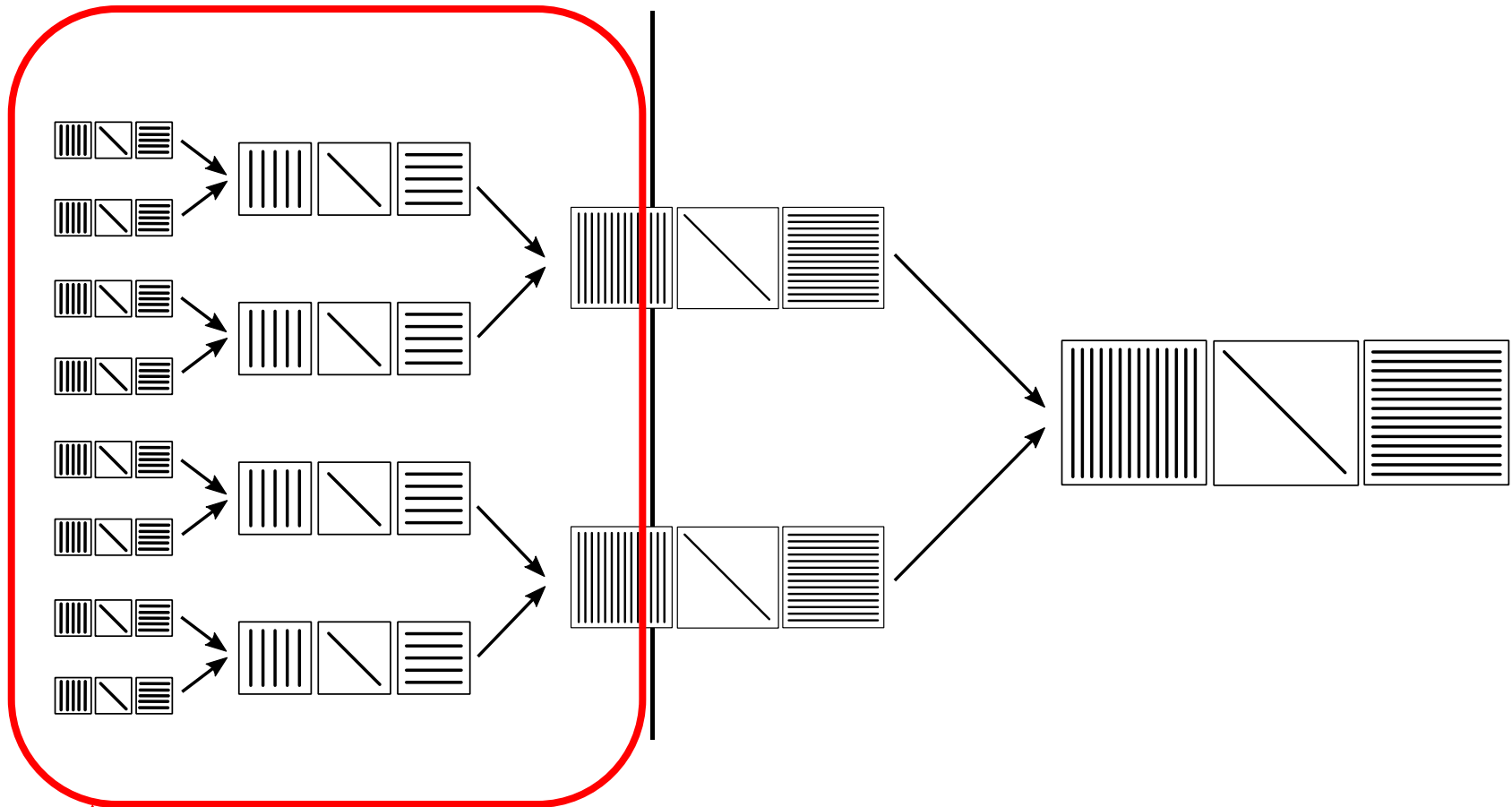
Allocate subproblems to both GPUs and CPUs

Model performance as a power function
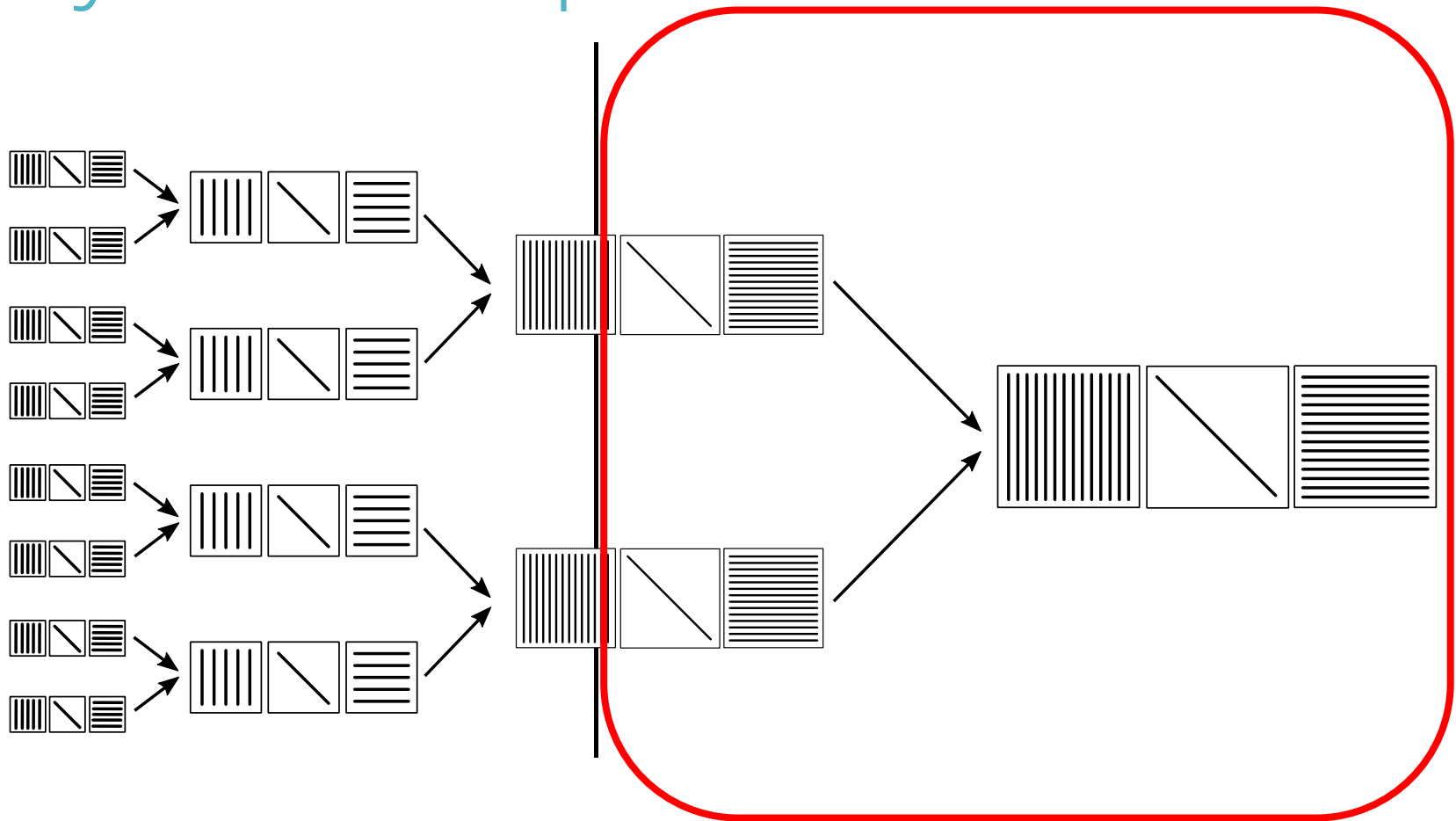
**Profiler** fits parameters using least-squares

performance

subproblem size

log(performance)

R² = 0.9566

log(subproblem size)

# Hybrid computation



Solve many subproblems in parallel

# Hybrid computation

Solve each subproblem by parts

# Results

**Per-GPU peak memory usage**
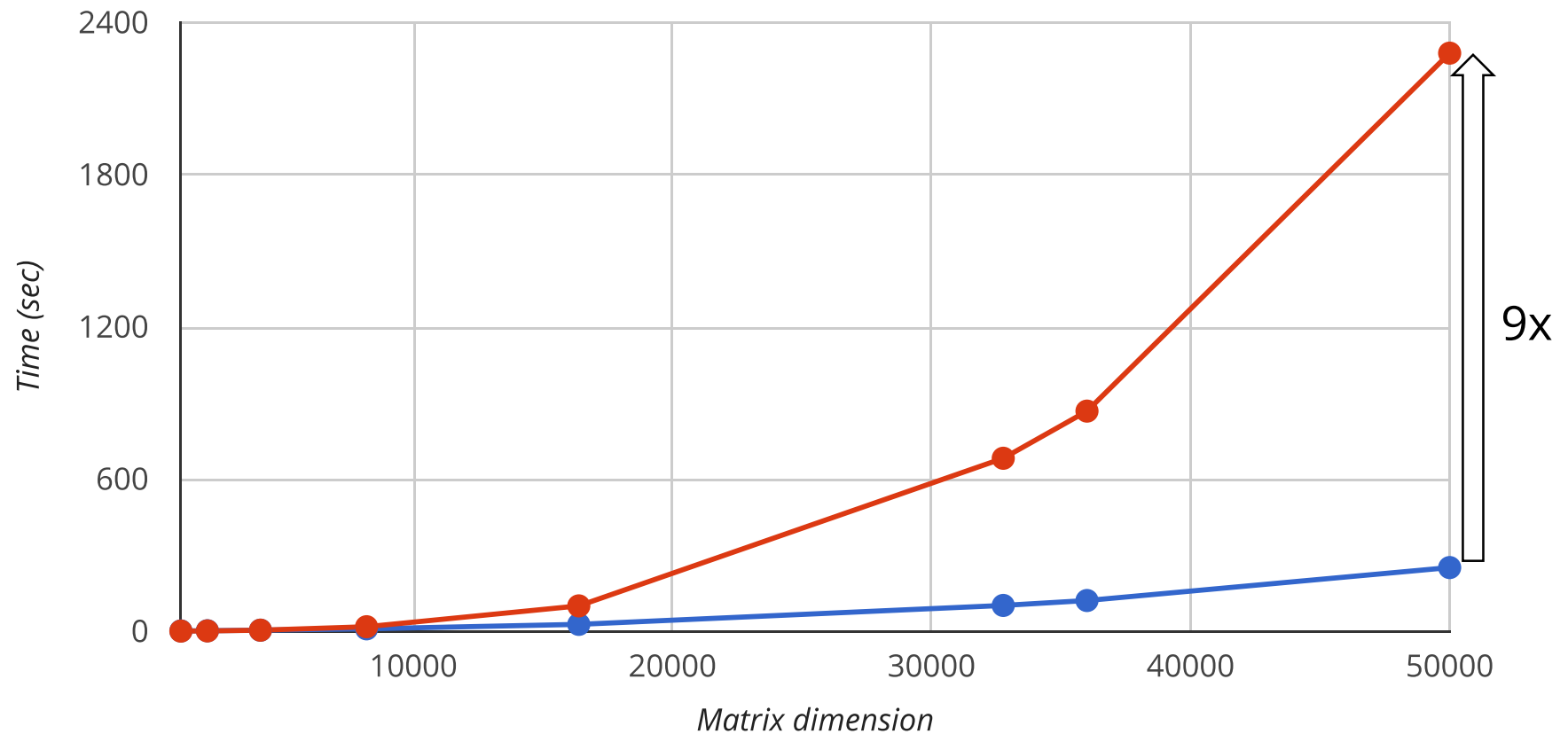
Scales to 50k * 50k matrix
With 4 GB of GPU memory

Resident memory size (MB)

4000

3000

2000

1000

0

0          10000          20000          30000          40000

*Matrix dimension*

Main memory: 64 GB
GPU memory: 5 GB per GPU

# Results

## Performance: vs. multicore CPU



Time (sec) — y-axis: 0, 600, 1200, 1800, 2400
Matrix dimension — x-axis: 10000, 20000, 30000, 40000, 50000

9x

CPU: dual Intel® Xeon® E5-2620
GPU: 4 Nvidia Tesla® K20c

■ GPUs +CPUs   ■ CPUs only

# Conclusion

Out-of-core approach overcomes memory limitation on the GPU

Hybrid computation with profiling delivers reasonable performance

# Acknowledgment

Trinity College, Student Research Program

Nvidia Corporation, CUDA Teaching Center Program

# Any questions?