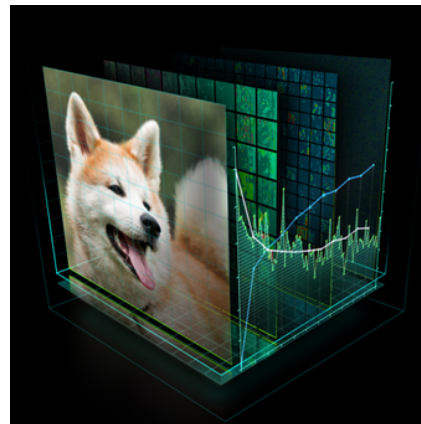# Introduction
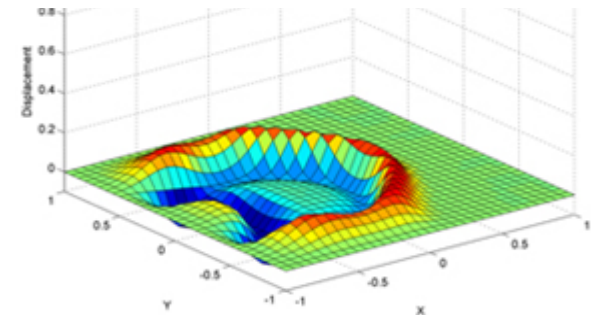
- Every new desktop/laptop is now equipped with a graphic processing unit (GPU).

- GPU = Massively Parallel Architecture.

- For most of their life, such GPUs are idle.

- General Purpose GPU applications:

Bioinformatics                Deep Learning                Numerical Analysis
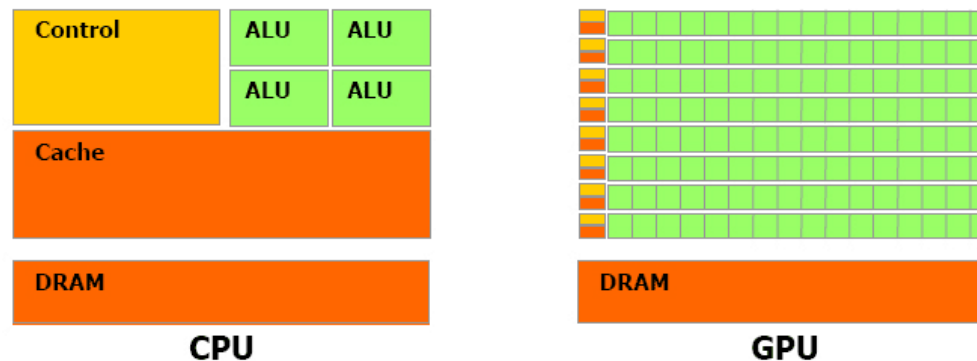                                                           MathWorks MATLAB

# (Distributed) Constraint Optimization

- A *(D)COP* is a tuple $<\textbf{X}, \textbf{D}, \textbf{F}, (A, \alpha)>$, where:

  - *X* is a set of variables.

  - *D* is a set of finite domains.

  - *F* is a set of utility functions:  $f_i : \times_{x_j \in scope(f_i)} D_j \mapsto \mathbb{N} \cup \{0, -\infty\}$

  - *A* is a set of agents, controlling the variables in *X*.

  - $\alpha$ maps variables to agents.

  - GOAL: Find a utility maximal assignment.

$$\mathbf{x}^* = \arg\max_{\mathbf{x}} \mathbf{F}(\mathbf{x})$$

$$= \arg\max_{\mathbf{x}} \sum_{f \in \mathbf{F}} f\left(\mathbf{x}\big|_{\mathrm{scope}(f)}\right)$$
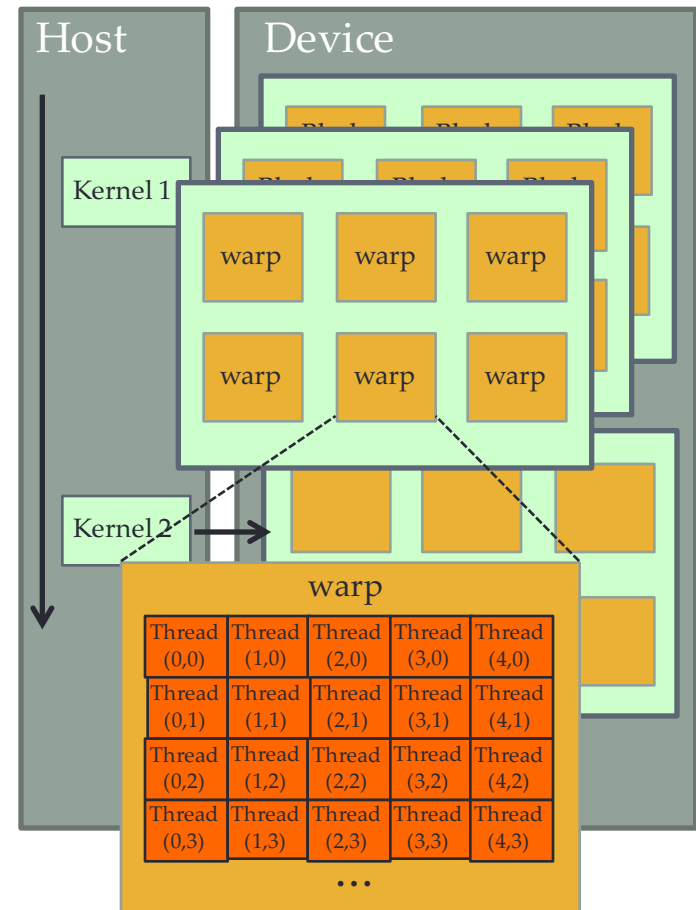
# Graphical Processing Units (GPUs)

- A GPU is a massive parallel architecture:
  - Thousands of multi-threaded **computing cores**.
  - Very high **memory bandwidths**.
  - ~80% of transistors devoted to data processing rather than caching.



- **However:**
  - GPU cores are slower than CPU cores.
  - GPU memories have different sizes and access times.
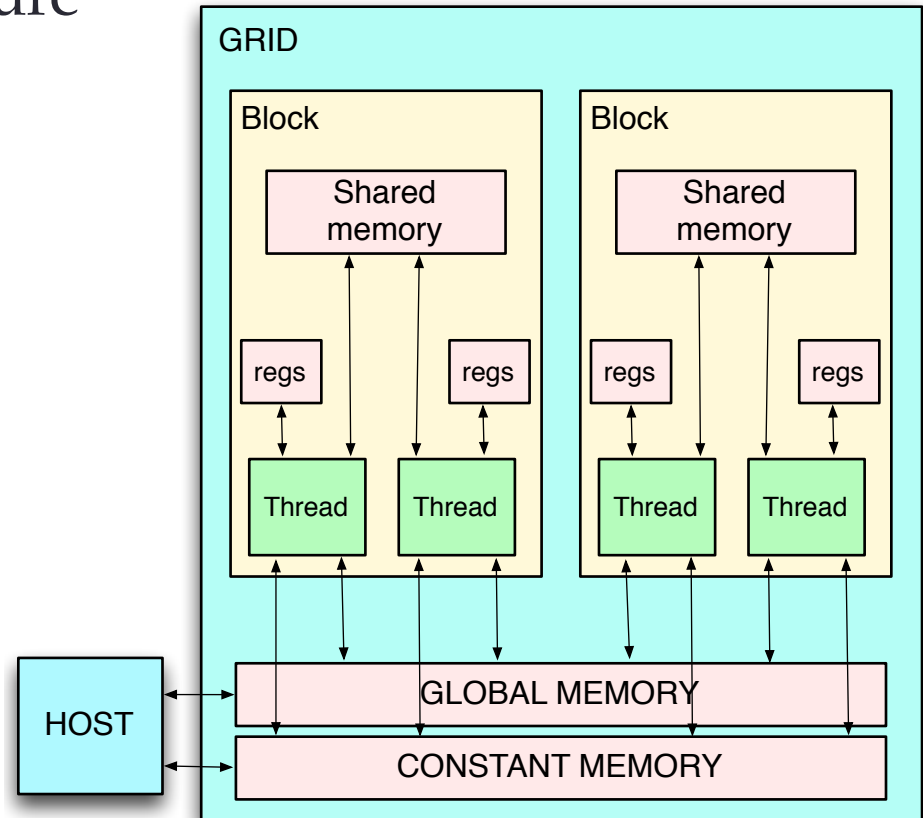  - GPU programming is more challenging and time consuming.

# Execution Model

- A Thread is the **basic parallel unit**.

- **Threads** are organized into a Block.
- Several **warps** are scheduled for the execution of a GPU function.

- Several Streaming Multiprocessors, (SD) scheduled in parallel.
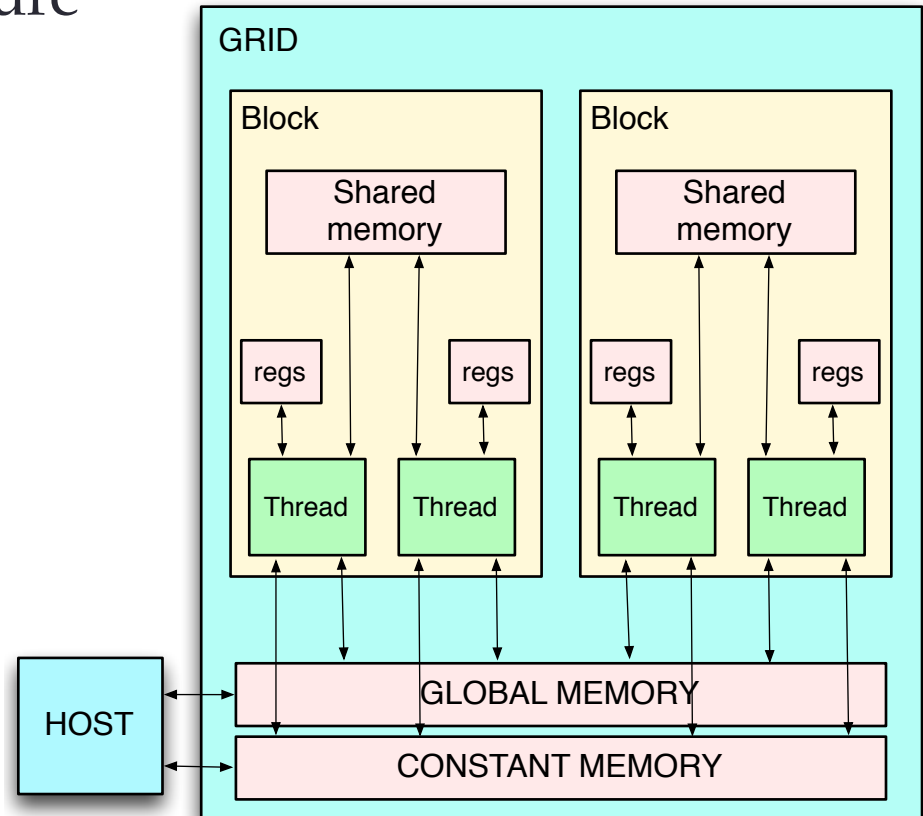- Single Instruction Multiple Thread (SIMT) parallel model.

# Memory Hierarchy

- The GPU memory architecture is rather involved.

- Registers

- Shared memory

- Global memory
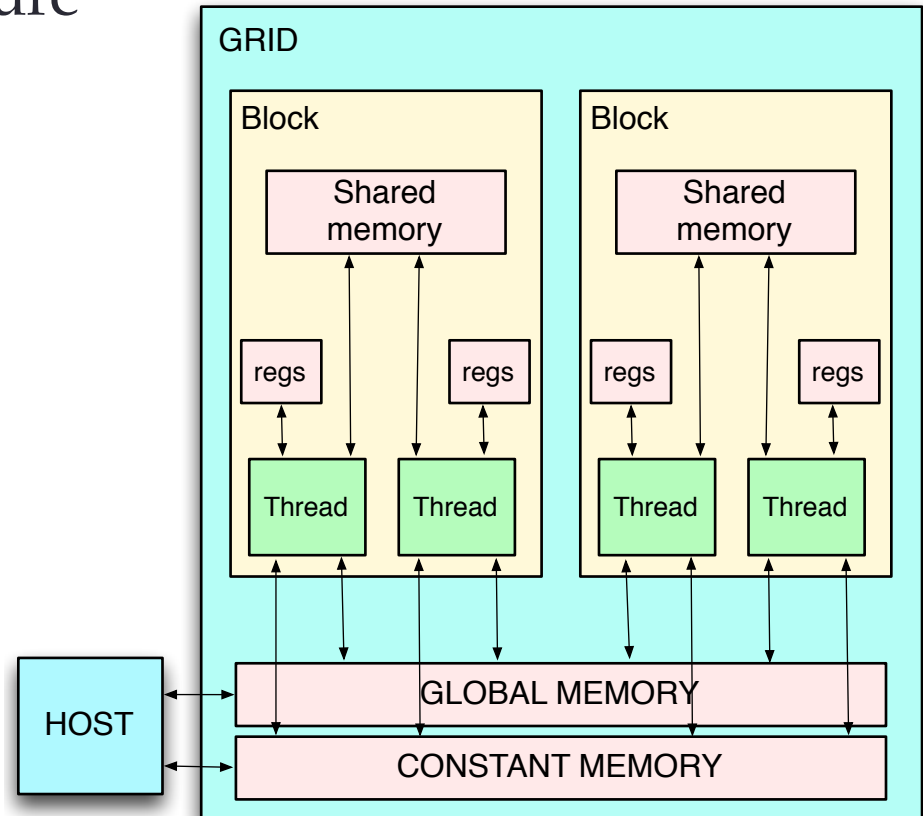
# Memory Hierarchy

- The GPU memory architecture is rather involved.

- Registers
  - **Fastest**;
  - Only accessible by a **thread**;
  - Lifetime of a thread.

- Shared memory

- Global memory

# Memory Hierarchy

- The GPU memory architecture is rather involved.

- Registers

- Shared memory
  - **Extremely fast**;
  - Highly parallel;
  - Restricted to a **block**.

- Global memory

# Memory Hierarchy

- The GPU memory architecture is rather involved.

- Registers

- Shared memory

- Global memory
  - Typically implemented in DRAM;
  - **High access latency** (400-800 cycles);
  - Potential of **traffic congestion**.

# Memory Hierarchy

- The GPU memory architecture is rather involved.

- Registers

- Shared memory

- Global memory

- **Challenge**: using memory effectively -- likely requires to redesign the algorithm.

# CUDA: Compute Unified Device Architecture

## Host

## Device

# CUDA: Compute Unified Device Architecture

**Host**                                    **Device**

cudaMalloc(&deviceV, sizeV);        data            **Global Memory**

cudaMemcpy(deviceV, hostV,
              sizeV, ...)

# CUDA: Compute Unified Device Architecture

**Host**                                                                **Device**

cudaKernel( )

cudaKernel<nThreads, nBlocks>( )     Kernel invocation     **Global Memory**

# CUDA: Compute Unified Device Architecture

**Host**                                                             **Device**



cudaMemcpy(hostV, deviceV,        ⟵ data        **Global Memory**
            sizeV, ...)

# Bucket Elimination and DPOP

- Dynamic Programming procedures to solve (D)COPs.

- Both procedures rely on the use of two operators:

- **Projection Operator:** $\pi_{-xi}(f_{ij})$

| $\mathbf{x_i}$ | $x_j$ | U |
|---|---|---|
| 0 | 0 | 5 |
| 0 | 1 | 8 |
| 1 | 0 | 20 |
| 1 | 1 | 3 |

$f_{ij}$

$\longrightarrow$

| $x_j$ | U |
|---|---|
| 0 | 20 |
| 1 | 8 |

- **Aggregation Operator:** $f_{ij} + f_{ik}$

# Bucket Elimination and DPOP

- Dynamic Programming procedures to solve (D)COPs.

- Both procedures rely on the use of two operators:

- **Projection Operator:** $\pi_{-xi}(f_{ij})$

| $\mathbf{x_i}$ | $x_j$ | U |
|------|------|-----|
| 0 | 0 | 5 |
| 0 | 1 | 8 |
| 1 | 0 | 20 |
| 1 | 1 | 3 |

$f_{ij}$

$\text{max}(5, \mathbf{20})$

| $x_j$ | U |
|------|-----|
| 0 | 20 |
| 1 | 8 |

- **Aggregation Operator:** $f_{ij} + f_{ik}$

# Bucket Elimination and DPOP

- Dynamic Programming procedures to solve (D)COPs.

- Both procedures rely on the use of two operators:

- **Projection Operator:** $\pi_{-xi}(f_{ij})$

| $\mathbf{x_i}$ | $x_j$ | U |
|---|---|---|
| 0 | 0 | 5 |
| 0 | 1 | 8 |
| 1 | 0 | 20 |
| 1 | 1 | 3 |

$f_{ij}$

max(**8**, 3)

| $x_j$ | U |
|---|---|
| 0 | 20 |
| 1 | 8 |

- **Aggregation Operator:** $f_{ij} + f_{ik}$

# Bucket Elimination and DPOP

- Dynamic Programming procedures to solve (D)COPs.

- Both procedures rely on the use of two operators:

- **Projection Operator:** $\pi_{-xi}(f_{ij})$

- **Aggregation Operator:** $f_{ij} + f_{ik}$

| $x_i$ | $x_j$ | U |
|-------|-------|----|
| 0 | 0 | 5 |
| 0 | 1 | 8 |
| 1 | 0 | 20 |
| 1 | 1 | 3 |

| $x_i$ | $x_k$ | U |
|-------|-------|----|
| 0 | 0 | 2 |
| 0 | 1 | 6 |
| 1 | 0 | 11 |
| 1 | 1 | 4 |

$$f_{ij}$$

| $x_i$ | $x_j$ | $x_k$ | U |
|-------|-------|-------|----|
| 0 | 0 | 0 | 7 |
| 0 | 0 | 1 | 11 |
| 0 | 1 | 0 | 10 |
| 0 | 1 | 1 | 14 |
| . . . | | | |

. . .

# Bucket Elimination and DPOP

- Dynamic Programming procedures to solve (D)COPs.
- Both procedures rely on the use of two operators:
- **Projection Operator:** $\pi_{-xi}(f_{ij})$
- **Aggregation Operator:** $f_{ij} + f_{ik}$

| $x_i$ | $x_j$ | U |
|-------|-------|-----|
| 0 | 0 | 5 |
| 0 | 1 | 8 |
| 1 | 0 | 20 |
| 1 | 1 | 3 |

$5 + 2 = \mathbf{7}$

| $x_i$ | $x_k$ | U |
|-------|-------|-----|
| 0 | 0 | 2 |
| 0 | 1 | 6 |
| 1 | 0 | 11 |
| 1 | 1 | 4 |

$f_{ij}$

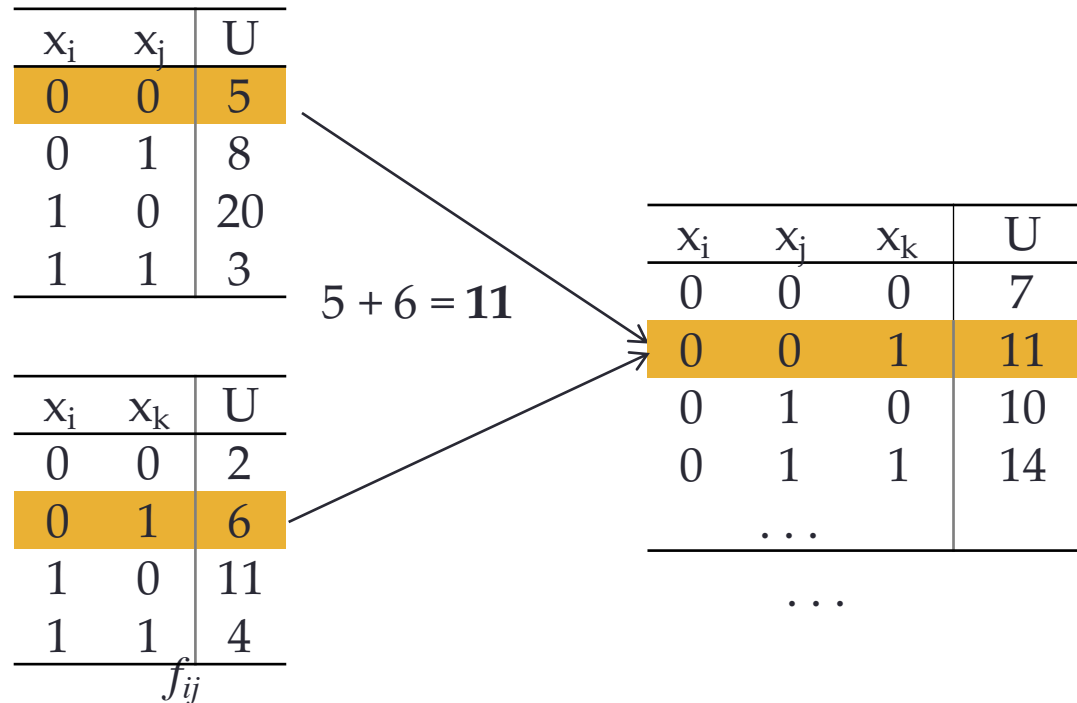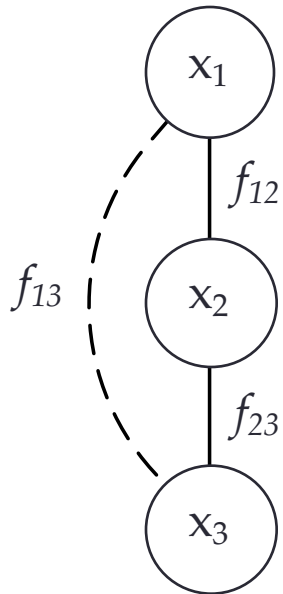| $x_i$ | $x_j$ | $x_k$ | U |
|-------|-------|-------|-----|
| 0 | 0 | 0 | 7 |
| 0 | 0 | 1 | 11 |
| 0 | 1 | 0 | 10 |
| 0 | 1 | 1 | 14 |
| . . . | | | |
| . . . | | | |

# Bucket Elimination and DPOP

- Dynamic Programming procedures to solve (D)COPs.
- Both procedures rely on the use of two operators:
- **Projection Operator:** $\pi_{-xi}(f_{ij})$
- **Aggregation Operator:** $f_{ij} + f_{ik}$

| $x_i$ | $x_j$ | U |
|-------|-------|-----|
| 0 | 0 | 5 |
| 0 | 1 | 8 |
| 1 | 0 | 20 |
| 1 | 1 | 3 |

$5 + 6 = \mathbf{11}$

| $x_i$ | $x_k$ | U |
|-------|-------|-----|
| 0 | 0 | 2 |
| 0 | 1 | 6 |
| 1 | 0 | 11 |
| 1 | 1 | 4 |

$f_{ij}$

| $x_i$ | $x_j$ | $x_k$ | U |
|-------|-------|-------|-----|
| 0 | 0 | 0 | 7 |
| 0 | 0 | 1 | 11 |
| 0 | 1 | 0 | 10 |
| 0 | 1 | 1 | 14 |
| . . . | | | |

. . .

# Bucket Elimination and DPOP

1.  Imposes an ordering on the problem's variables.
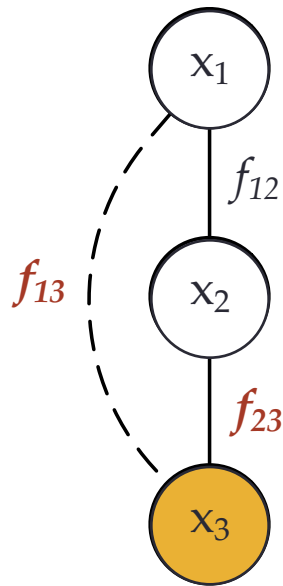


$X = \{x_1, x_2, x_3\}$

$F = \{f_{12}, f_{13}, f_{23}\}$

# Bucket Elimination and DPOP

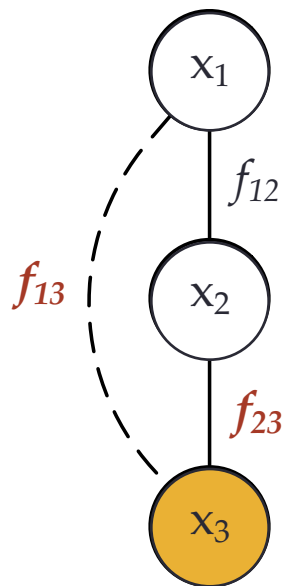2. Selects the variable $x_i$ with highest priority, and it creates a **bucket:**



$$B_i = \left\{ f_j \in \mathbf{F} \mid x_i \in scope(f_j) \land \right.$$
$$\left. i = \max\{k \mid x_k \in scope(f_j)\} \right\}$$

$X = \{x_1, x_2, x_3\}$     $B_3 = \{f_{13}, f_{23}\}$

$F = \{f_{12}, f_{13}, f_{23}\}$

# Bucket Elimination and DPOP

3. It computes a new utility function $f_i'$ by aggregating the functions in $B_i$ and projecting out $x_i$



| $f_{13}$ | $x_1$ | $x_3$ | U |
|---|---|---|---|
| | 0 | 0 | 5 |
| | 0 | 1 | 8 |
| | 1 | 0 | 20 |
| | 1 | 1 | 3 |

| $f_{23}$ | $x_2$ | $x_3$ | U |
|---|---|---|---|
| | 0 | 0 | 5 |
| | 0 | 1 | 8 |
| | 1 | 0 | 20 |
| | 1 | 1 | 3 |

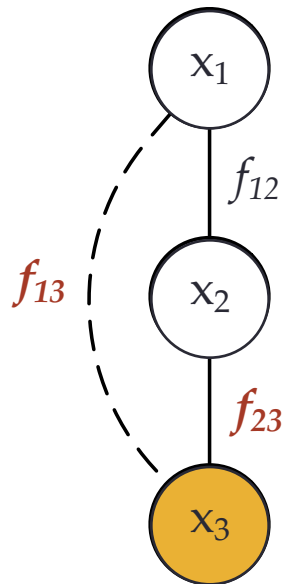| $f_3'$ | $x_1$ | $x_2$ | U | $(x_3)$ |
|---|---|---|---|---|
| | 0 | 0 | **max**(5 + 5, 8 + 8) = **16** | 1 |

$$f_3' = \pi_{-x3} \left( f_{13} + f_{23} \right)$$

$X = \{x_1, x_2, x_3\}$     $B_3 = \{f_{13}, f_{23}\}$

$F = \{f_{12}, f_{13}, f_{23}\}$

# Bucket Elimination and DPOP

3. It computes a new utility function $f_i'$ by aggregating the functions in $B_i$ and projecting out $x_i$



| $f_{13}$ | $x_1$ | $x_3$ | U |
|---|---|---|---|
| | 0 | 0 | 5 |
| | 0 | 1 | 8 |
| | 1 | 0 | 20 |
| | 1 | 1 | 3 |

| $f_{23}$ | $x_2$ | $x_3$ | U |
|---|---|---|---|
| | 0 | 0 | 5 |
| | 0 | 1 | 8 |
| | 1 | 0 | 20 |
| | 1 | 1 | 3 |

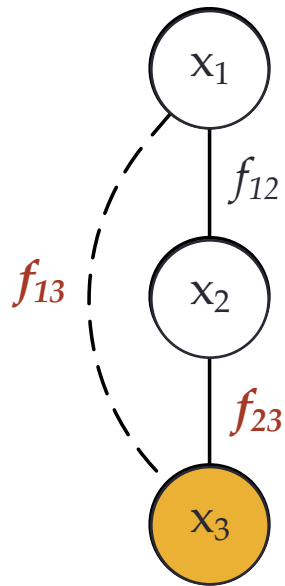| $f_3'$ | $x_1$ | $x_2$ | U | $(x_3)$ |
|---|---|---|---|---|
| | 0 | 0 | 16 | 1 |
| | 0 | 1 | **max**$(5 + 20, 8 + 3) =$ **25** | 0 |

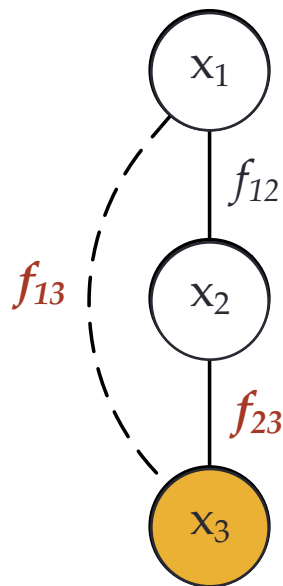$$f_3' = \pi_{-x_3}\left(f_{13} + f_{23}\right)$$

$$X = \{x_1, x_2, x_3\} \qquad B_3 = \{f_{13}, f_{23}\}$$

$$F = \{f_{12}, f_{13}, f_{23}\}$$

# Bucket Elimination and DPOP

3. It computes a new utility function $f_i'$ by aggregating the functions in $B_i$ and projecting out $x_i$



| $f_{13}$ | $x_1$ | $x_3$ | U |
|---|---|---|---|
| | 0 | 0 | 5 |
| | 0 | 1 | 8 |
| | 1 | 0 | 20 |
| | 1 | 1 | 3 |

| $f_{23}$ | $x_2$ | $x_3$ | U |
|---|---|---|---|
| | 0 | 0 | 5 |
| | 0 | 1 | 8 |
| | 1 | 0 | 20 |
| | 1 | 1 | 3 |

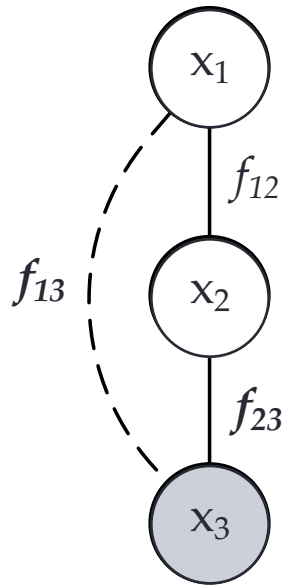| $f_3'$ | $x_1$ | $x_2$ | U | $(x_3)$ |
|---|---|---|---|---|
| | 0 | 0 | 16 | 1 |
| | 0 | 1 | 25 | 0 |
| | 1 | 0 | **max**$(20 + 5, 3 + 8) = $ **25** | 0 |

$$f_3' = \pi_{-x3} (f_{13} + f_{23})$$

$X = \{x_1, x_2, x_3\}$    $B_3 = \{f_{13}, f_{23}\}$

$F = \{f_{12}, f_{13}, f_{23}\}$

# Bucket Elimination and DPOP

3. It computes a new utility function $f_i'$ by aggregating the functions in $B_i$ and projecting out $x_i$



| $f_{13}$ | $x_1$ | $x_3$ | U |
|---|---|---|---|
| | 0 | 0 | 5 |
| | 0 | 1 | 8 |
| | 1 | 0 | 20 |
| | 1 | 1 | 3 |

| $f_{23}$ | $x_2$ | $x_3$ | U |
|---|---|---|---|
| | 0 | 0 | 5 |
| | 0 | 1 | 8 |
| | 1 | 0 | 20 |
| | 1 | 1 | 3 |

| $f_3'$ | $x_1$ | $x_2$ | U | $(x_3)$ |
|---|---|---|---|---|
| | 0 | 0 | 16 | 1 |
| | 0 | 1 | 25 | 0 |
| | 1 | 0 | 25 | 0 |
| | 1 | 0 | **max**$(20 + 20, 3 + 3) = $ **40** | 0 |

$$f_3' = \pi_{-x3} \left( f_{13} + f_{23} \right)$$

$X = \{x_1, x_2, x_3\}$     $B_3 = \{f_{13}, f_{23}\}$

$F = \{f_{12}, f_{13}, f_{23}\}$

# Bucket Elimination and DPOP

4.  It updates the set of variables: $\mathbf{X} \leftarrow \mathbf{X} \setminus \{x_i\}$
5.  It updates the set of functions: $\mathbf{F} \leftarrow (\mathbf{F} \cup \{f_i'\}) \setminus B_i$
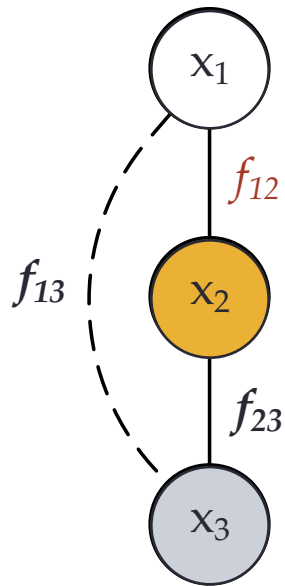


$X = \{x_1, x_2\}$        $B_2 = \{f_{13}, f_3'\}$

$F = \{f_{12}, f_3'\}$

# Bucket Elimination and DPOP

Repeat...
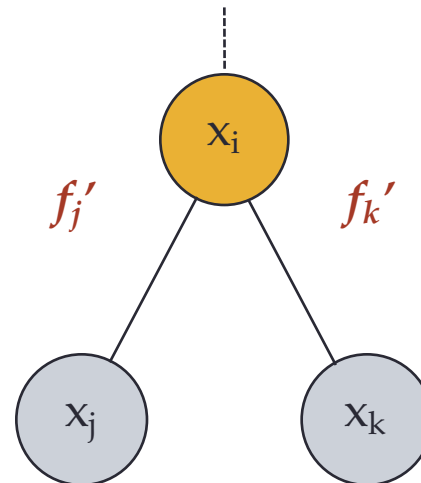
- **DPOP** is a distributed version of BE.
- It operates on a Pseudotree ordering of the constraint graph.

$X = \{x_1, x_2\}$

$F = \{f_{12}, f_3{}'\}$

# GPU-(D)BE

- BE and DPOP complexity: $O(d^{w^*})$.
  d = max. domain size; w* = induced width of the constraint graph.
- Can the projection and aggregator operators be executed in parallel?
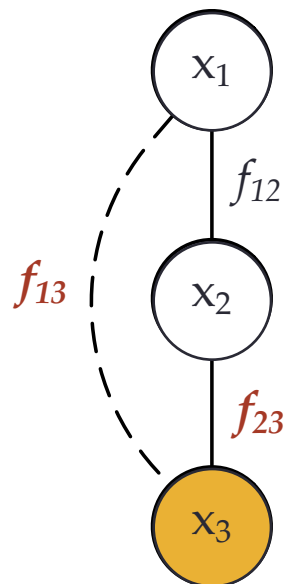- Do they fit the SIMT parallel model?

| $x_1$ | $x_3$ | U |
|---|---|---|
| 0 | 0 | 5 |
| 0 | 1 | 8 |
| 1 | 0 | 20 |
| 1 | 1 | 3 |

| $x_2$ | $x_3$ | U |
|---|---|---|
| 0 | 0 | 5 |
| 0 | 1 | 8 |
| 1 | 0 | 20 |
| 1 | 1 | 3 |

| $x_1$ | $x_2$ | U |
|---|---|---|
| 0 | 0 | **max**(5 + 5, 8 + 8) = **16** |
| 0 | 1 | **max**(5 + 20, 8 + 3) = **25** |
| 1 | 0 | **max**(20 + 5, 3 + 8) = **25** |
| 1 | 0 | **max**(20 + 20, 3 + 3) = **40** |

$$f_3' = \pi_{-x3}\left(f_{13} + f_{23}\right)$$

# GPU-(D)BE

- BE and DPOP complexity: $O(d^{w*})$.
  d = max. domain size; w* = induced width of the constraint graph.

- Can the projection and aggregator operators be executed in parallel?

- Do they fit the SIMT parallel model?

- **Obs.**: The computation of each row of the Utility tables is independent from the computation of other rows.



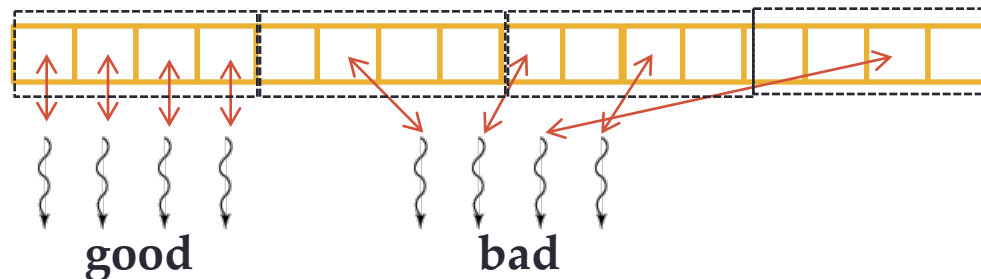| $x_1$ | $x_3$ | U |
|---|---|---|
| 0 | 0 | 5 |
| 0 | 1 | 8 |
| 1 | 0 | 20 |
| 1 | 1 | 3 |

| $x_2$ | $x_3$ | U |
|---|---|---|
| 0 | 0 | 5 |
| 0 | 1 | 8 |
| 1 | 0 | 20 |
| 1 | 1 | 3 |

| $x_1$ | $x_2$ | U |
|---|---|---|
| 0 | 0 | **max**(5 + 5, 8 + 8) = **16** |
| 0 | 1 | **max**(5 + 20, 8 + 3) = **25** |
| 1 | 0 | **max**(20 + 5, 3 + 8) = **25** |
| 1 | 0 | **max**(20 + 20, 3 + 3) = **40** |

$$f_3{'} = \pi_{-x3}\left(f_{13} + f_{23}\right)$$
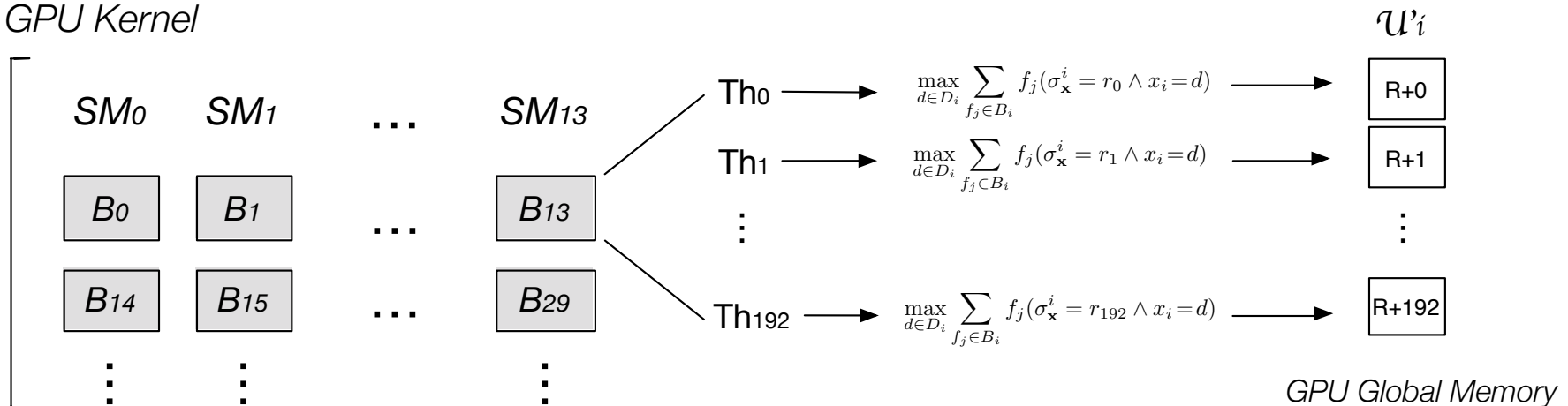
# Algorithm design and data structure

- Limit the amount of host-device **data transfers**.
  - Static Entities: require a **single** data transaction.
    - Variables; Domains; Utility functions; Constraint Graph.
  - Dynamic Entities: might require **multiple** data transactions.
    - Utility tables.

- Minimize the accesses to the **global memory**.

  - Padding Utility Tables' rows; Perfect hashing.

- Ensure data accesses are **coalesced**.

  - Mono-dimensional array organization;



**good**          **bad**

# Parallel Projection and Aggregation

- Mapping between the $f_i'$ table rows and the *CUDA* blocks:
  - Each **thread** in a block is associated to the computations of one permutation of values in *scope($f_i'$)*.
  - 1 **block** = *64k* threads *(1 ≤ k ≤ 16)*.
    - k depends on the architecture and it is chosen so to maximize the number of threads that can be scheduled concurrently.
- **Obs.**: Max number of parallel $f_i'$ table rows is *M=|SM|64k*
  - In our experiments, *|SMs| = 14* and *k = 3*. Thus *M = 2688*.

*GPU Kernel*



$$\max_{d \in D_i} \sum_{f_j \in B_i} f_j(\sigma_{\mathbf{x}}^i = r_0 \wedge x_i = d)$$

$$\max_{d \in D_i} \sum_{f_j \in B_i} f_j(\sigma_{\mathbf{x}}^i = r_1 \wedge x_i = d)$$

$$\max_{d \in D_i} \sum_{f_j \in B_i} f_j(\sigma_{\mathbf{x}}^i = r_{192} \wedge x_i = d)$$
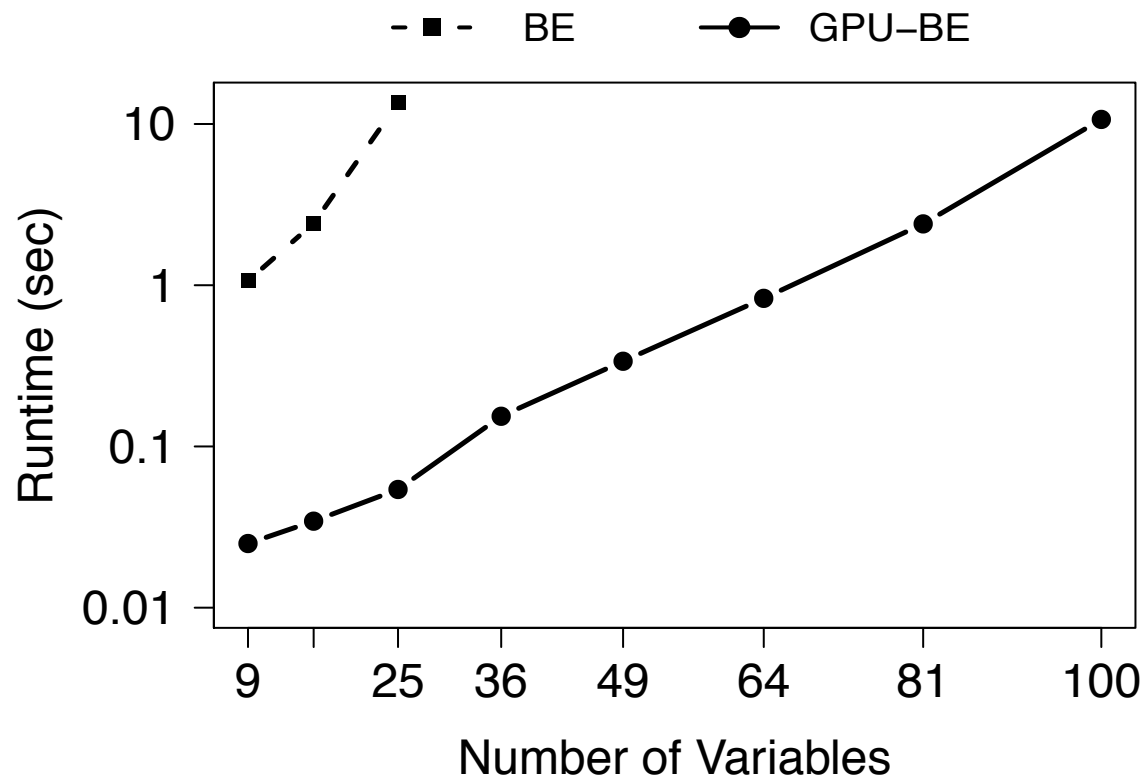
*GPU Global Memory*

# Experimental Results

- |Di| = 5;
- p2 = 0.5;
- timeout = 300s

**Regular Grid Networks**

- CPU: 2.3GHz, 128 GB RAM
- GPU: 14 SMs, 837MHz.



**Speedup**
avg. max.: 125.1x
avg. min.:  42.6x

- Similar trends at increasing |$D_i$|.
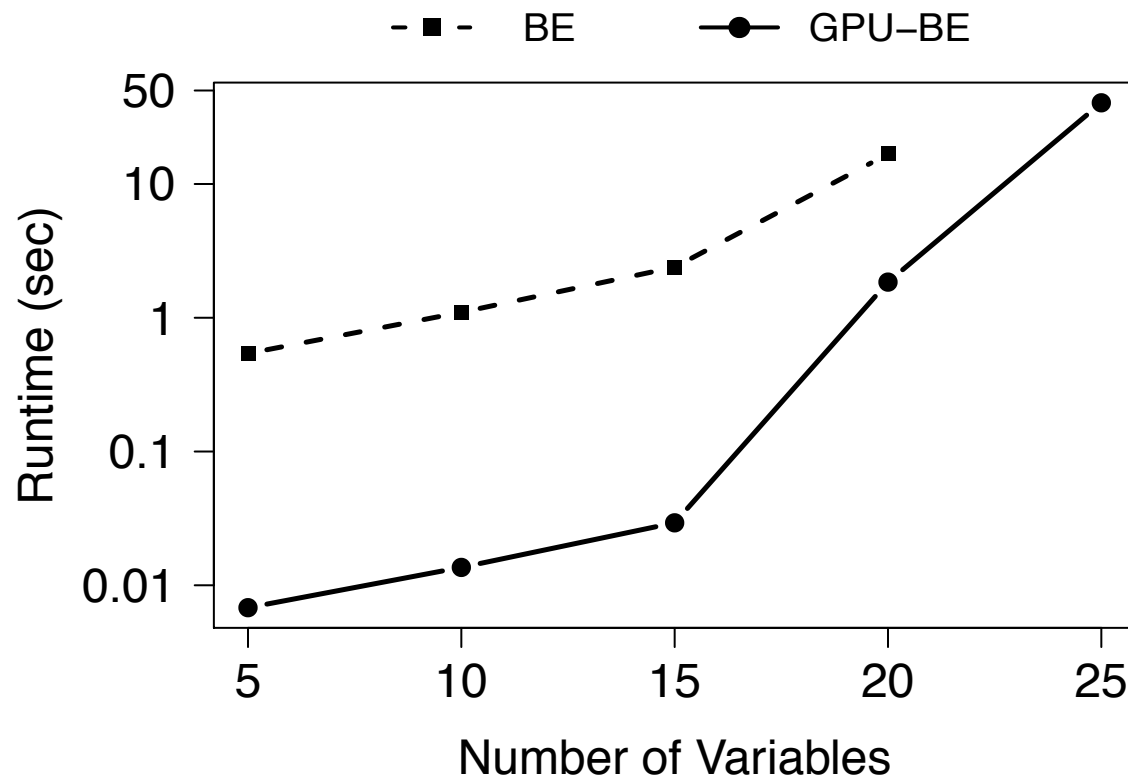- Similar trends for DPOP vs GPU-DBE

# Experimental Results

- |Di| = 5;
- p2 = 0.5; p1 = 0.3;
- timeout = 300s

**Random Networks**

- CPU: 2.3GHz, 128 GB RAM
- GPU: 14 SMs, 837MHz.
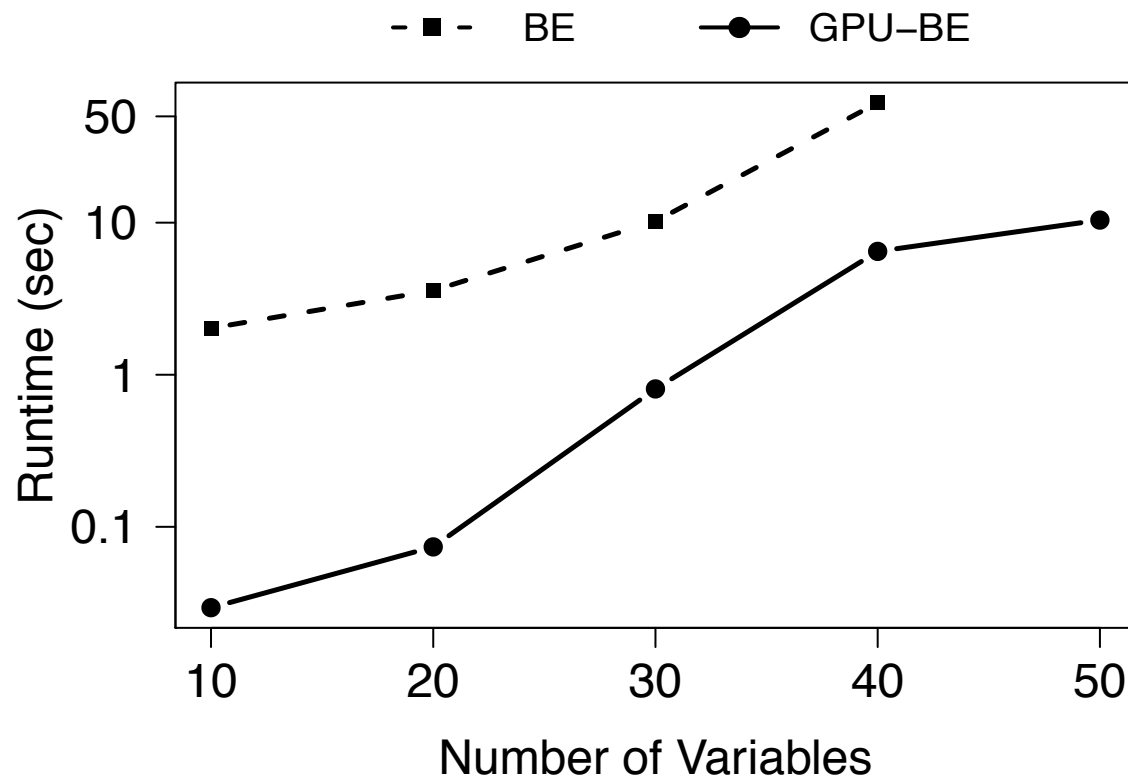


**Speedup**

avg. max.:  69.3.x
avg. min.:   16.1x

- Similar trends at increasing $|D_i|$.
- Similar trends for DPOP vs GPU-DBE

# Experimental Results

- |Di| = 5;
- p2 = 0.5;
- timeout = 300s

- CPU: 2.3GHz, 128 GB RAM
- GPU: 14 SMs, 837MHz.

**Scale Free Networks**



**Speedup**
avg. max.:  34.9.x
avg. min.:   9.5x

- Similar trends at increasing $|D_i|$.
- Similar trends for DPOP vs GPU-DBE

# Lesson Learned #1
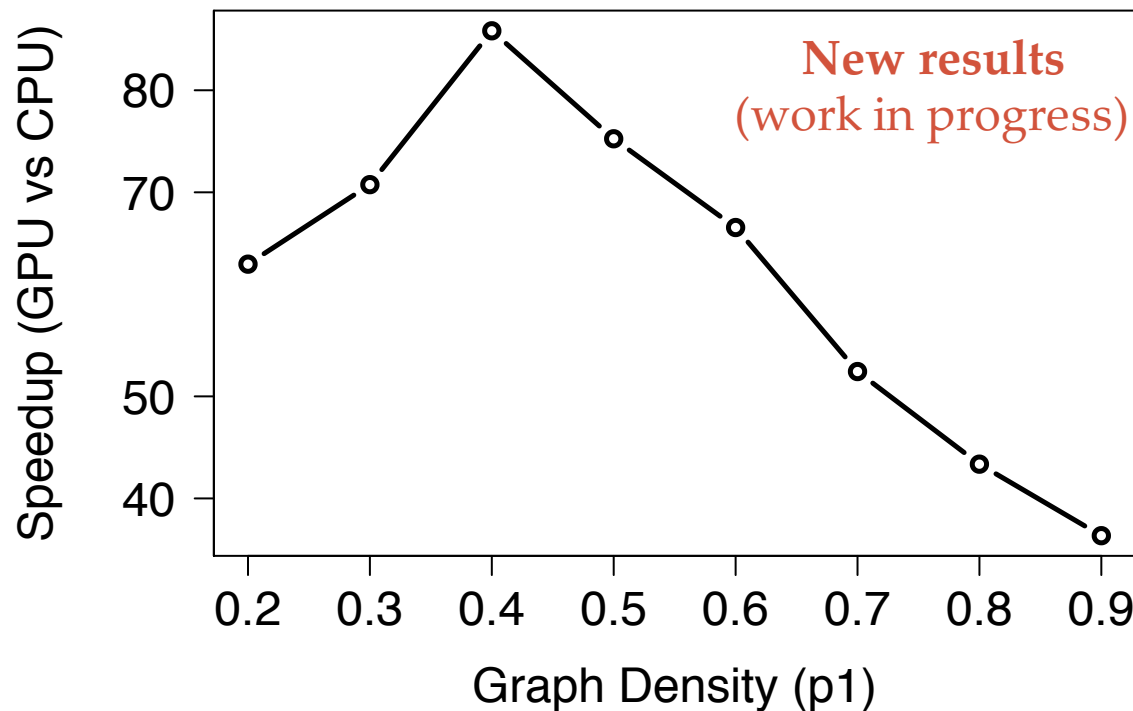
- The $f_i'$ table size increases exponentially with $w^*$.

- Limited GPU global memory (2GB).

- $f_i'$ table + $B_i$ tables, to be used in the aggregation operation, might exceed global memory capacity!

- Partition $f_i'$ computations in multiple chunks.

- Alternates GPU and CPU to compute $f_i'$.
  - GPU: Aggregates the functions in $B_i$ excluding those which do not fit in the global memory.
  - CPU:  Aggregates the other functions in $B_i$;
          Projects out the variable $x_i$.

# Experimental Results

- $|A| = 10$; $|D_i| = 5$;
- $p_2 = 0.5$;
- timeout = 300s

**Random Networks**

- CPU: 2.3GHz, 128 GB RAM
- GPU: 14 SMs, 837MHz.

**New results**
**(work in progress)**

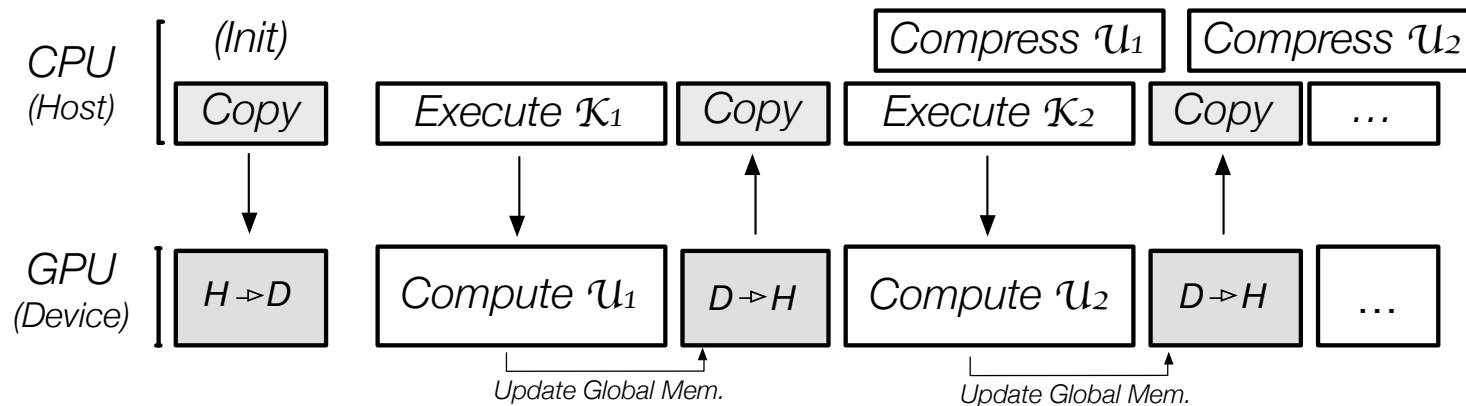**Phase Transition**
$p_1 = 0.4$

- small $p_1$ correspond to smaller w*

# Lesson Learned #2

- Host and Device **concurrency**.

- Possible when the $f_i'$ tables are computed in **chunks**.

- It may hide host-device data transfers as byproduct.

# Discussion

- Exploiting the integration of CPU and GPU is a key factor to obtain competitive solver performance.

- How to determine good tradeoffs of such integration?
- **GPU**:
  - Repeated, non memory intensive operations;
  - Operations requiring regular memory access;
- **CPU:**
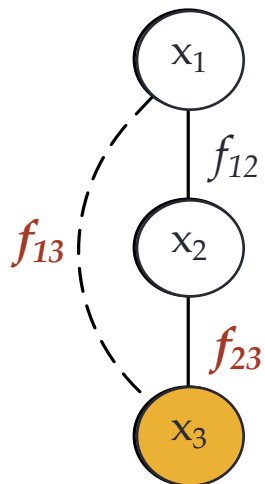  - Memory intensive operations;

# Conclusions

- Exploit GPU-style parallelism from DP-based (D)COPs resolution methods.

- GPU-(D)BE: Exploits GPUs to parallelizes the **aggregation** and **projection** operators.

- Observed different speedup, ranging from 34.9 to 125.1, based on several network topologies.

- Discussed several possible optimization techniques.

- **FUTURE WORK**:
  - Exploit GPUs in DP-based propagators.
  - Investigate GPUs in higher form of consistency.

# Exploiting GPUs in Solving (Distributed) Constraint Optimization Problems with Dynamic Programming

F. Fioretto, T. Le, E. Pontelli, W. Yeoh, T. Son

# Thank You!



Ferdinando Fioretto

New Mexico State University,
University of Udine

Email: ffiorett@cs.nmsu.edu
Web: www.cs.nmsu.edu/~ffiorett