# GPU Programming

## Farshad Khunjush

Department of Computer Science and Engineering

Shiraz University

Fall 2025

# Optimization Techniques

Farshad Khunjush

# Performance Considerations
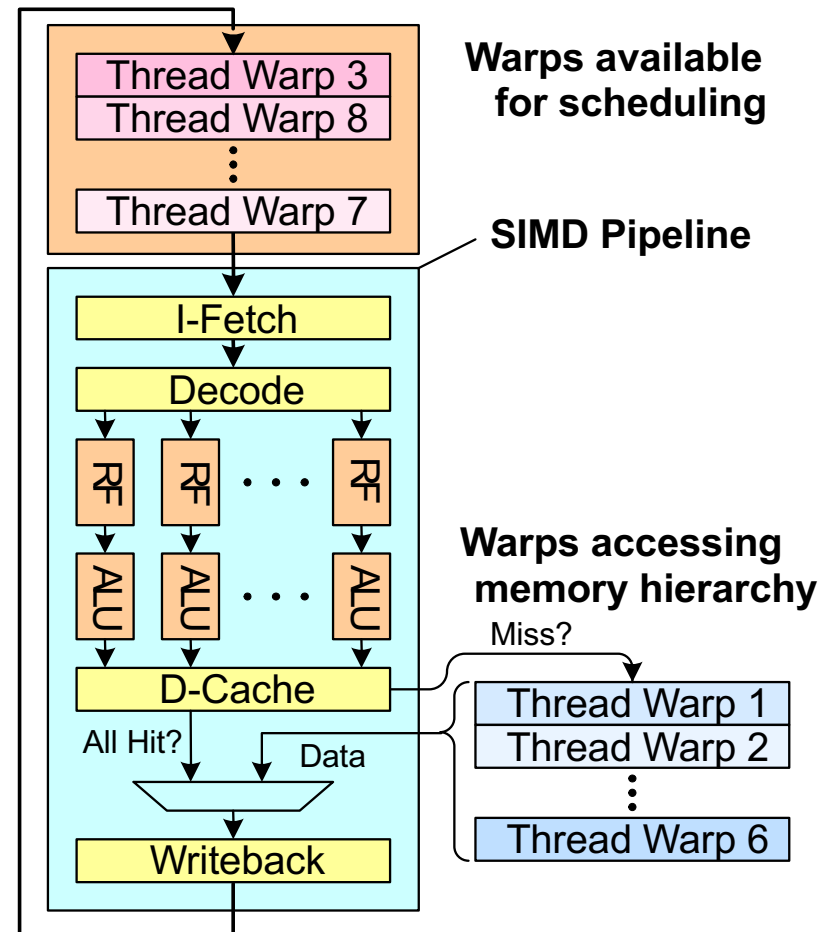
- Main bottlenecks
  - Global memory access
  - CPU-GPU data transfers
- Memory access
  - Latency hiding
    - Occupancy
  - Memory coalescing
  - Data reuse
    - Shared memory usage
- SIMD (Warp) Utilization: Divergence
- Other considerations
  - Atomic operations: Serialization
  - Data transfers between CPU and GPU
    - Overlap of communication and computation
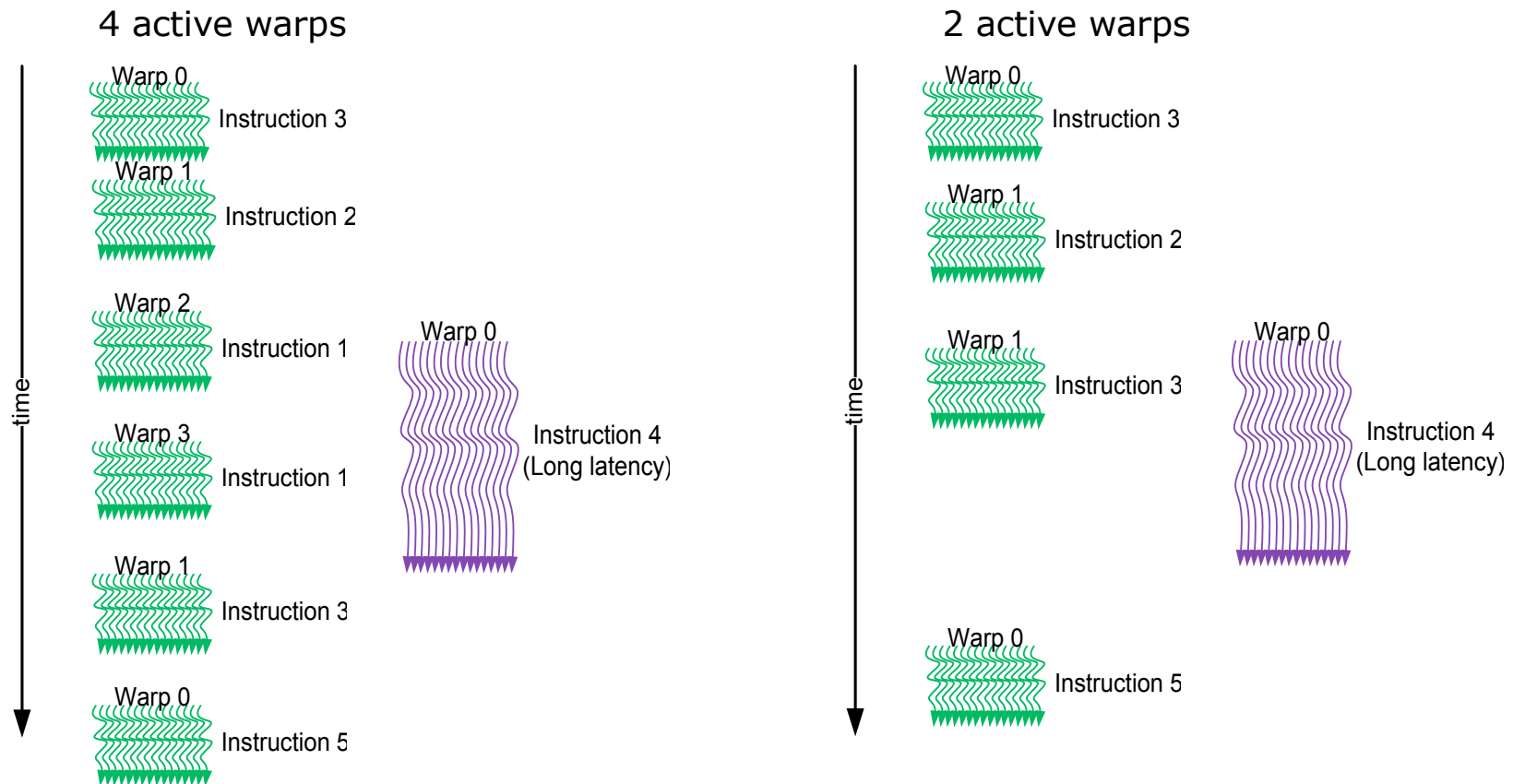
# Memory Access

# Latency Hiding via Warp-Level FGMT

- Warp: A set of threads that execute the same instruction (on different data elements)

- Fine-grained multithreading
  - One instruction per thread in pipeline at a time (No interlocking)
  - Interleave warp execution to hide latencies
- Register values of all threads stay in register file
- FGMT enables long latency tolerance
  - Millions of pixels

**Warps available for scheduling**

Thread Warp 3
Thread Warp 8
⋮
Thread Warp 7

**SIMD Pipeline**

I-Fetch

Decode

RF   RF   · · ·   RF

ALU   ALU   · · ·   ALU

D-Cache

All Hit?          Data

Writeback

**Warps accessing memory hierarchy**

Miss?

Thread Warp 1
Thread Warp 2
⋮
Thread Warp 6

5

Slide credit: Tor Aamodt

# Latency Hiding and Occupancy

- FGMT can hide long latency operations (e.g., memory accesses)
- Occupancy: ratio of active warps to the maximum number of warps per GPU core

4 active warps

Warp 0
Instruction 3

Warp 1
Instruction 2

Warp 2
Instruction 1

Warp 3
Instruction 1

Warp 1
Instruction 3

Warp 0
Instruction 5

time

Warp 0
Instruction 4
(Long latency)

2 active warps

Warp 0
Instruction 3

Warp 1
Instruction 2

Warp 1
Instruction 3

Warp 0
Instruction 5

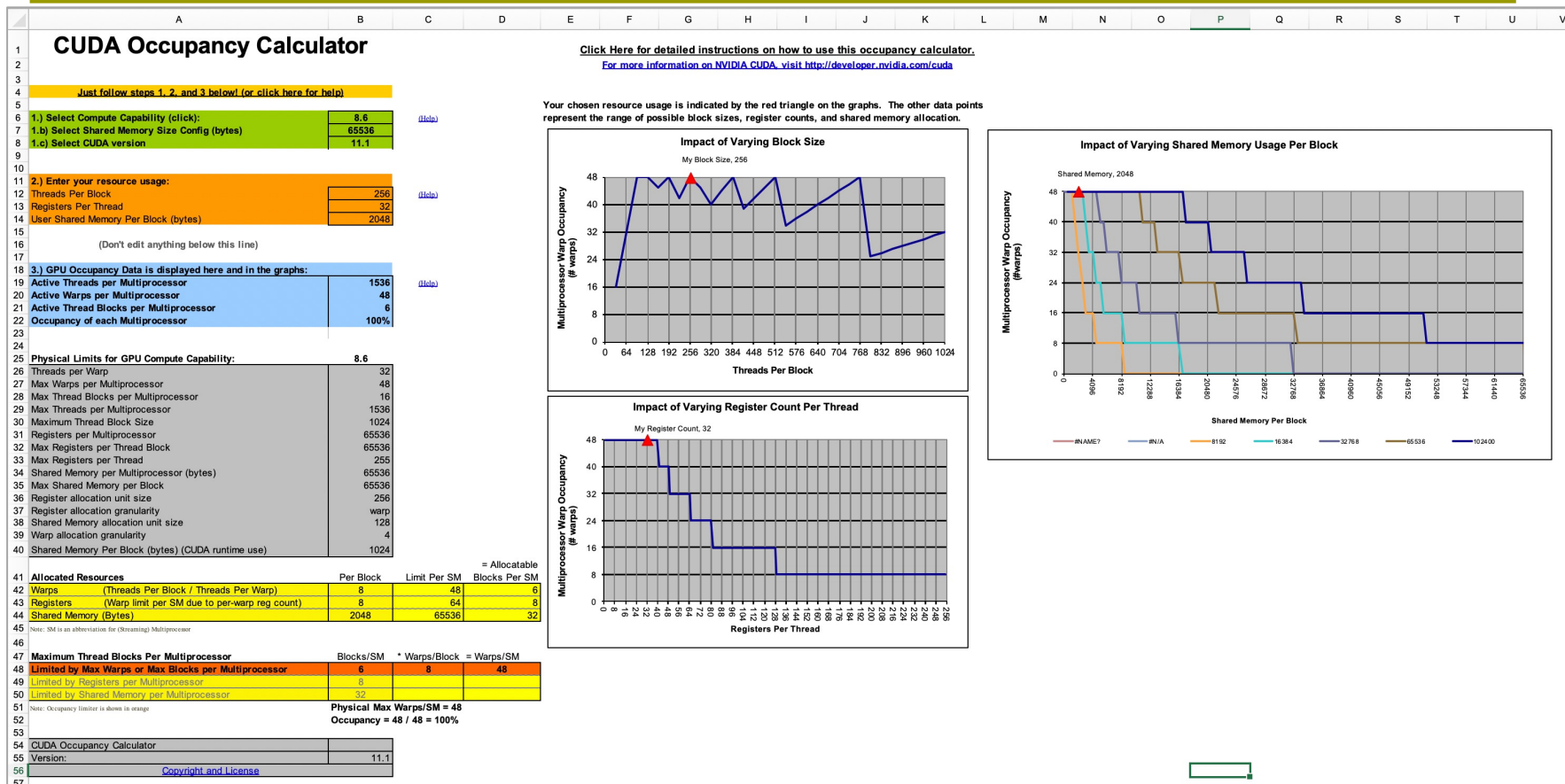time

Warp 0
Instruction 4
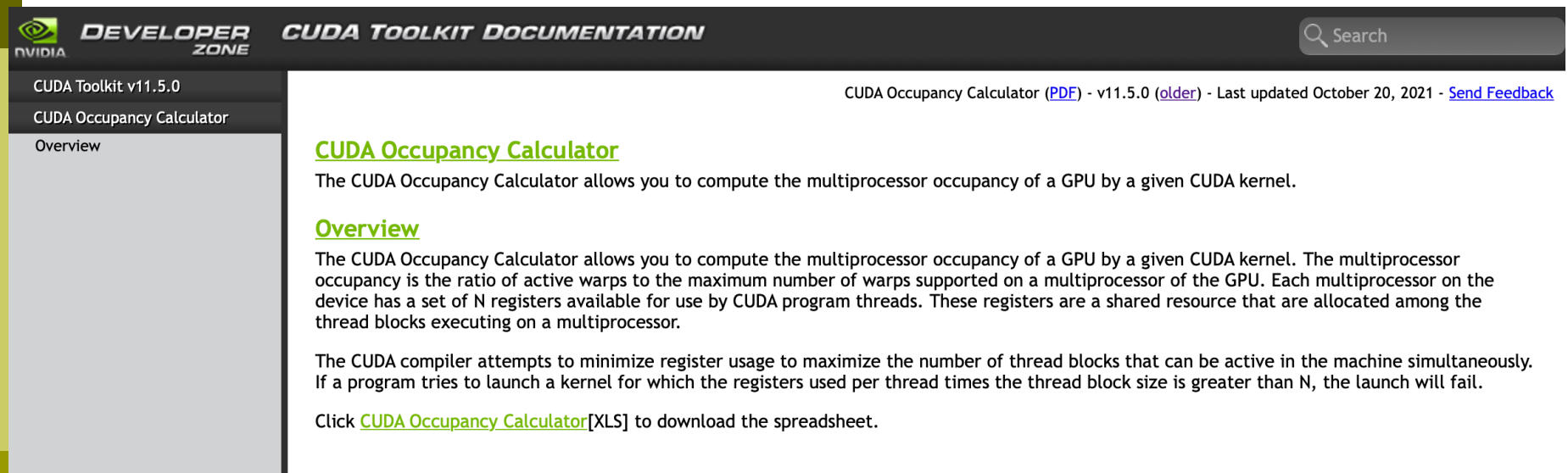(Long latency)

# Occupancy

- GPU core, a.k.a. SM, resources (typical values)
  - Maximum number of warps per SM (64)
  - Maximum number of blocks per SM (32)
  - Register usage (256KB)
  - Shared memory usage (64KB)

- Occupancy calculation
  - Number of threads per block (defined by the programmer)
  - Registers per thread (known at compile time)
  - Shared memory per block (defined by the programmer)

# CUDA Occupancy Calculator (I)



https://docs.nvidia.com/cuda/cuda-occupancy-calculator/CUDA_Occupancy_Calculator.xls

# CUDA Occupancy Calculator (II)



**CUDA Occupancy Calculator**

The CUDA Occupancy Calculator allows you to compute the multiprocessor occupancy of a GPU by a given CUDA kernel.

**Overview**

The CUDA Occupancy Calculator allows you to compute the multiprocessor occupancy of a GPU by a given CUDA kernel. The multiprocessor occupancy is the ratio of active warps to the maximum number of warps supported on a multiprocessor of the GPU. Each multiprocessor on the device has a set of N registers available for use by CUDA program threads. These registers are a shared resource that are allocated among the thread blocks executing on a multiprocessor.

The CUDA compiler attempts to minimize register usage to maximize the number of thread blocks that can be active in the machine simultaneously. If a program tries to launch a kernel for which the registers used per thread times the thread block size is greater than N, the launch will fail.

Click CUDA Occupancy Calculator[XLS] to download the spreadsheet.

https://docs.nvidia.com/cuda/cuda-occupancy-calculator/CUDA_Occupancy_Calculator.xls

# Optimization: Algorithm

- **Maximize parallelism**
  - Avoid thread sync. if possible
  - Minimize control-flow divergence

- **Maximize use of available memory bandwidth**
  - GPUs have four types of memory (more later)

- **Maximize arithmetic intensity (compute/access)**
  - GPU spends its transistors on ALUs, not memory
    - Re-compute rather than cache
  - Avoid CPU-GPU data transfer:
    - Do more computation on GPU
    - Or, compute on CPU if low parallelism

# GPU Memory System

- Global and local memories
- Shared memory and L1 cache
- Registers
- Constant & Texture memories ◆ L2 cache

# GPU's Heterogeneous Memory Structures

- Designed originally for graphics computing
  - Specialized storage for various graphics data
  - Global, Constant, Texture
- Caches
  - For bandwidth improvement
  - Not latency (unlike CPU's)
  - Not always coherent
- Memory hierarchy
  - Localized connectivity (improve bandwidth)

# GPU Memory Hierarchy: Global Memory

- **Global memory** (per application):
  - Shared by all threads
  - GPU's main memory (separate HW from GPU core)
  - ~10GB, ~300 GB/s of BW and latency of ~400 cycles
  - Inter-grid communication

# GPU Memory Hierarchy: Local Memory

- ❑ <span style="color:red">Local memory</span> (per thread): **Thread**

  - ■ Private per thread

    **Local Memory**

  - ■ Everything on the stack that can't fit in registers
    - ❑ <span style="color:red">Register spilling</span>
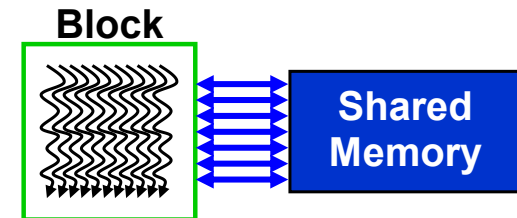
  - ■ Stored in global memory
    - ❑ Same <span style="color:red">latency</span> as global memory
    - ❑ Much <span style="color:red">slower</span> than registers!

# GPU Memory Hierarchy: Shared Memory

- <span style="color:red">Shared memory (per block):</span>
    - Shared within a thread block

    - **<span style="color:red">Managed by the programmer</span>**

    - Very fast, located in the SM
        - Latency: ~5ns
        - Bandwidth: ~1 TB/s

    - Same HW as L1 cache (64KB)
        - 16/32/48KB of L1 cache

    - Inter-thread communication

**Block**

**Shared Memory**

# GPU Memory Hierarchy: Caches

- **L1 Cache**:
    - Each SM has its own L1 cache
    - Older gen.: caches local & global memory
    - Newer gen.: only cache local memory
    - Same HW as shared memory
        - Configurable size: 16/32/48KB

- **L2 Cache**:
    - Shared by all SM's
    - Caches all global & local memory accesses
    - ~1MB size, ~500 GB/s of BW

# GPU Memory Hierarchy: Register File

❑ Registers:

- Stack variables declared in kernels
- Fastest access to data
  - ~10X faster than shared memory
- Bandwidth: a few 10s of TB/s
- Fundamental challenge in GPU microarchitecture

- Example (Fermi architecture):
  - 32K 32-bit registers per SM
  - 48 warps per SM
    - 1536 threads per SM → 21 registers/thread
  - 2MB register file
    - 16 SMs→128KB per SM

# GPU Memory Hierarchy: Constant & Texture

Historical leftover from graphics:

- Two more types of memory
  - But are not used that often
  - They are beneficial only for very specific types of apps

- Constant & Texture Memories:
  - Global with a special cache
  - Must be set from host before running kernel
    - Read-only over the course of a kernel execution
  - Can be used to reduce pressure on global memory

# Optimization: Exploit Shared Memory

- Hundreds of times faster than global memory

- Use shared memory as a <span style="color:red">managed scratchpad</span>
  - Bring data in from global memory
  - Operate in there (reuse)
  - Write results back to global memory

- Shared memory is fast as long as there are no bank conflicts

# Shared Memory Bank Conflicts

- Shared memory is organized in <span style="color:red">32 banks</span>
  - Each bank can service one address at a time
  - At most 32 simultaneous accesses

- <span style="color:red">Multiple simultaneous</span> accesses to the same bank
  - Different 4-byte words:
    - Bank conflict! - Conflicting accesses are serialized

  - The same 4-byte word:
    - Multicast – 1 fetch (could be different bytes within the word)

# Shared Memory Bank Conflicts

**No Bank Conflicts**
- Linear addressing stride == 1

| Thread 0 | → | Bank 0 |
| Thread 1 | → | Bank 1 |
| Thread 2 | → | Bank 2 |
| Thread 3 | → | Bank 3 |
| Thread 4 | → | Bank 4 |
| Thread 5 | → | Bank 5 |
| Thread 6 | → | Bank 6 |
| Thread 7 | → | Bank 7 |
| Thread 31 | → | Bank 31 |

**No Bank Conflicts**
- Random 1:1 Permutation

Thread 0, Thread 1, Thread 2, Thread 3, Thread 4, Thread 5, Thread 6, Thread 7, Thread 31 → Bank 0, Bank 1, Bank 2, Bank 3, Bank 4, Bank 5, Bank 6, Bank 7, Bank 31
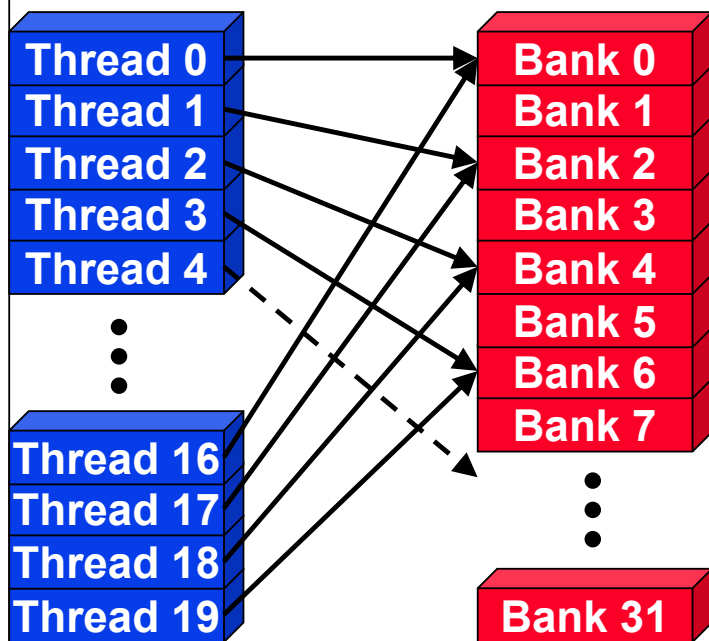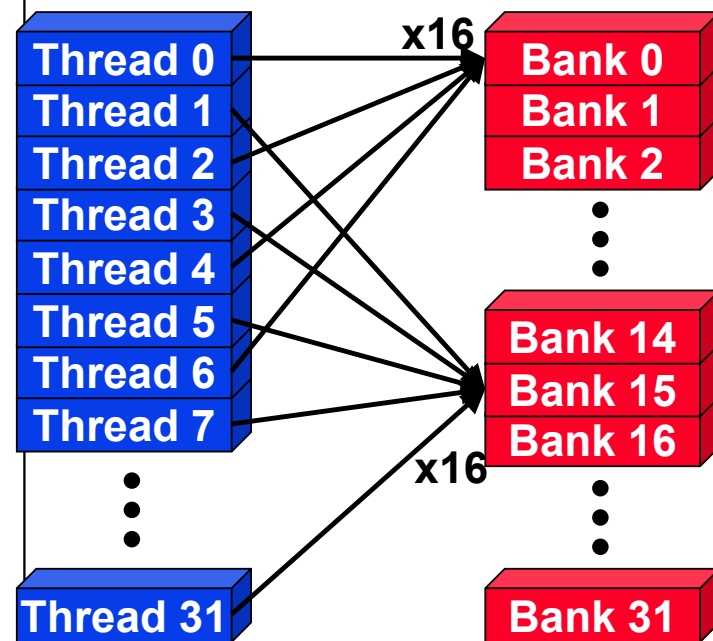
# Shared Memory Bank Conflicts

**2-way Bank Conflicts**

- Linear addressing stride == 2

**16-way Bank Conflicts**

- Linear addressing stride == 16

# How Addresses are Mapped to Banks

- Each bank has a BW of 32 bits per clock cycle
- Successive 32-bit words are assigned to successive banks
  - This is called memory interleaving

- Modern GPUs have 32 banks
  - So bank = address % 32

# Shared Memory Bank Conflicts - Example

- **Operating on 2D array in shared memory**
  - E.g., image processing

- **Example: 32x32 block**
  - Each thread processes a row
  - Threads in a block access the elements of a column simultaneously
    - E.g., Column 1 in purple
  - All 32 elements map to same bank
    - Elements are 32-bits apart in memory
    - 32-way bank conflict
    - 32 serialized accesses→SLOW!

**Bank Indices**

| 0 | 1 | 2 | 3 | 4 | 5 | ••• | 31 |
|---|---|---|---|---|---|-----|----|
| 0 | 1 | 2 | 3 | 4 | 5 | ••• | 31 |
| 0 | 1 | 2 | 3 | 4 | 5 | ••• | 31 |
| 0 | 1 | 2 | 3 | 4 | 5 | ••• | 31 |
| 0 | 1 | 2 | 3 | 4 | 5 | ••• | 31 |
| 0 | 1 | 2 | 3 | 4 | 5 | ••• | 31 |
| 0 | 1 | 2 | 3 | 4 | 5 | ••• | 31 |
| 0 | 1 | 2 | 3 | 4 | 5 | ••• | 31 |
| 0 | 1 | 2 | 3 | 4 | 5 | ••• | 31 |

# Shared Memory Bank Conflicts - Example

- **Solution 1: pad the rows**
  - Add one element to the end of each row
  - Now each thread goes to a different bank

**Bank Indices with Padding**

| 0 | 1 | 2 | 3 | 4 | 5 | ••• | 31 | 0 |
|---|---|---|---|---|---|-----|----|---|
| 1 | 2 | 3 | 4 | 5 | 6 | ••• | 0 | 1 |
| 2 | 3 | 4 | 5 | 6 | 7 | ••• | 1 | 2 |
| 3 | 4 | 5 | 6 | 7 | 8 | ••• | 2 | 3 |
| 4 | 5 | 6 | 7 | 8 | 9 | ••• | 3 | 4 |
| 5 | 6 | 7 | 8 | 9 | 10 | ••• | 4 | 5 |
| 6 | 7 | 8 | 9 | 10 | 11 | ••• | 5 | 6 |
| 7 | 8 | 9 | 10 | 11 | 12 | ••• | 7 | 8 |

| 31 | 0 | 1 | 2 | 3 | 4 | ••• | 30 | 31 |
|----|---|---|---|---|---|-----|----|----|

- **Solution 2: transpose first**
  - Suffer bank conflicts during transpose
  - But possibly save later conflicts
    - When passing multiple times over the same data

- **Solution 3: operate on columns instead of rows**
  - No bank conflict: one thread per bank
    Bank Indices with Padding

# Control-Flow Divergence

- **Threads in a warp might execute different paths**
  - Branch instructions in the code
- **Execute one path at a time**
  - Diverging threads will be disabled
  - Limits parallelism and lowers performance



**50% Performance Loss**

# Control-Flow Divergence

- A common case:
  - Avoid divergence when branch condition is a function of threadId (or tid)

- Example:

```
if (threadIdx.x > 2) {. . .}
```

  - Two different control paths for threads in a warp

# Control-Flow Divergence

- Example that may look divergent but it's not:

  ```
  if (threadIdx.x / WARP_SIZE > 2) {. . .}
  ```

  - Two different control paths, different warps, same TB
    - All threads in a warp follow the same path
    - Divergent across warps is OK!