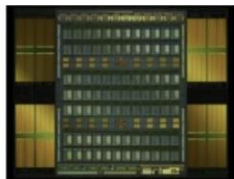


# GPU Architectures & Tensor Streaming Processors

Yağız Çakmak

# GPUs (Graphics Processing Units)

# Evolution of Recent GPUs (I)



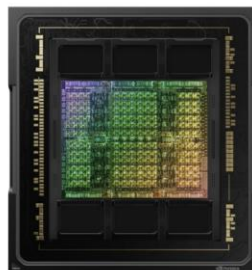
## Volta

>21 billion transistors  
815mm<sup>2</sup>  
TSMC 12nm FFN



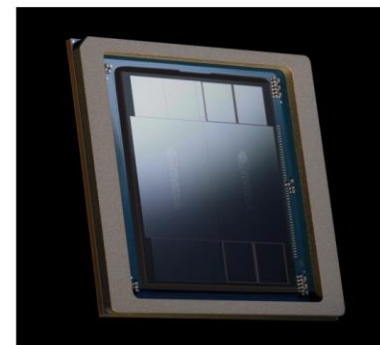
## Ampere

>54 billion transistors  
826 mm<sup>2</sup>  
TSMC N7



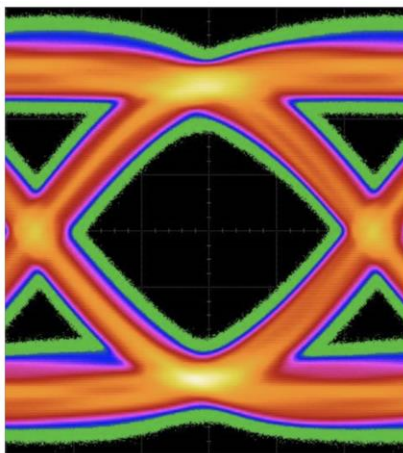
## Hopper

>80 billion transistors  
814 mm<sup>2</sup>  
TSMC 4N

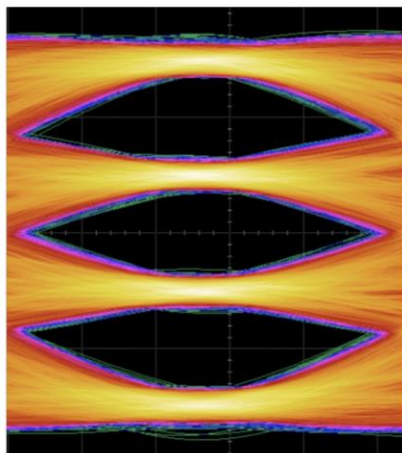


## Blackwell

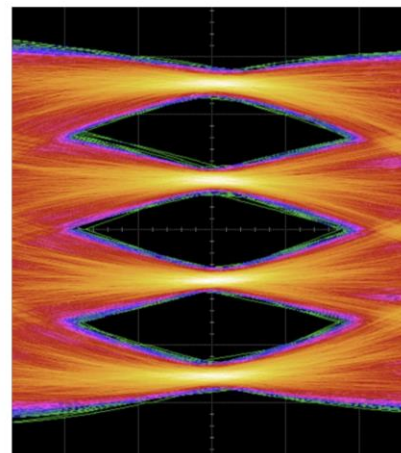
>208 billion transistors  
>1600 mm<sup>2</sup>  
TSMC 4NP



**Ampere | NVLink3**  
12 NVLinks | 50GB/s each  
x4@50Gbps-NRZ  
600GB/s total



**Hopper | NVLink4**  
18 NVLinks | 50GB/s each  
x2@100Gbps-PAM4  
900GB/s total

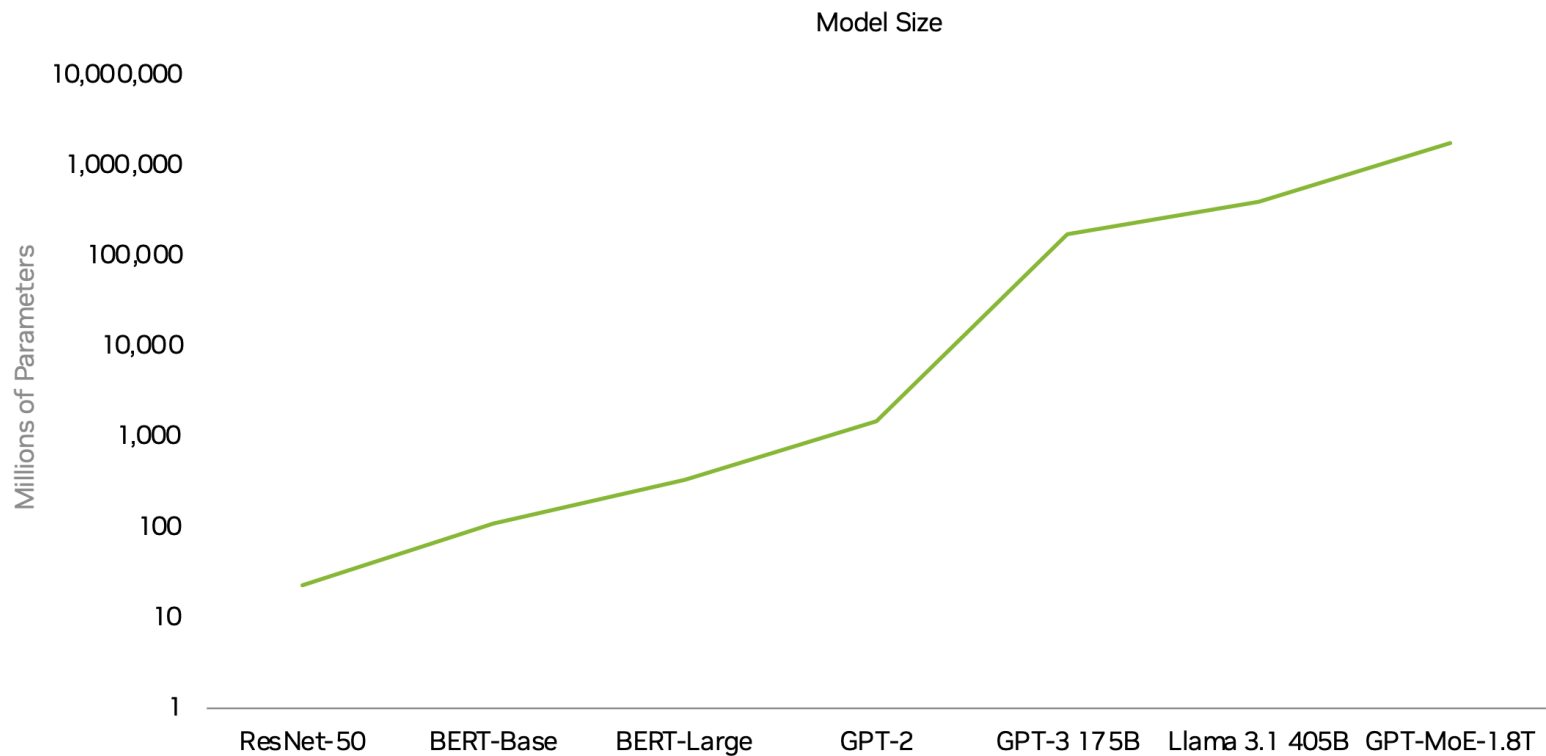


**Blackwell | NVLink5**  
18 NVLinks | 100GB/s each  
x2@200Gbps-PAM4  
1800GB/s total

# Multiple GPUs to Tackle Large Workloads

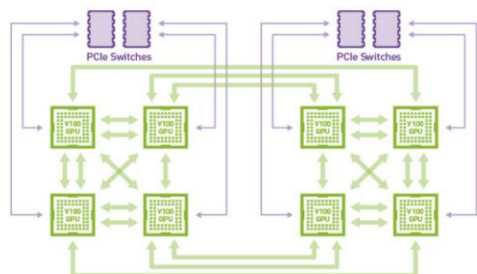
## AI Models Growing Exponentially

Need for multi-GPU inference at scale



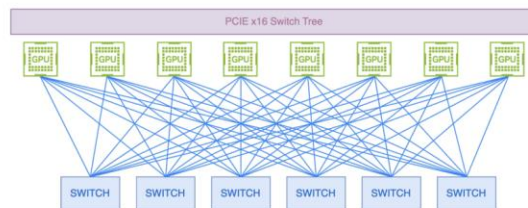
**New Capabilities | Trillions of Parameters | 70,000X Growth in a Decade**

# Evolution of Recent GPUs (II)



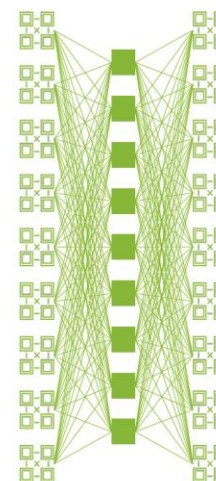
**2016**

Hybrid Cube Mesh NVLink technology



**2022**

3<sup>rd</sup> Gen NVLink Switch  
All-to-all connection among NVLink domain of 8 GPU



**2024**

4<sup>th</sup> Gen NVLink Switch Chip  
All-to-all connection among NVLink domain of 72 GPU

# GPUs are SIMD Engines Underneath

---

- The **instruction pipeline operates like a SIMD pipeline** (e.g., an array processor)
- However, the **programming is done using threads**, NOT SIMD instructions
- To understand this, let's go back to our parallelizable code example
- But, before that, let's distinguish between
  - **Programming Model (Software)**
  - vs.
  - **Execution Model (Hardware)**

# Programming Model vs. Hardware Execution Model

---

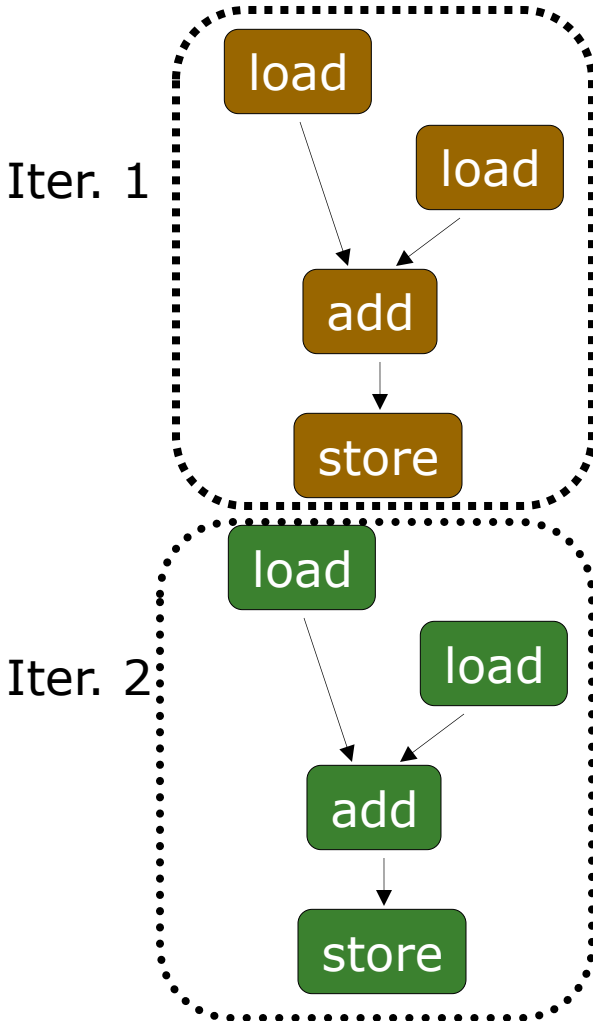
- Programming Model refers to **how the programmer expresses the code**
  - E.g., Sequential (von Neumann), Data Parallel (SIMD), Dataflow, Multi-threaded (MIMD, SPMD), ...
- Execution Model refers to **how the hardware executes the code underneath**
  - E.g., Out-of-order execution, Vector processor, Array processor, Dataflow processor, Multiprocessor, Multithreaded processor, ...
- **Execution Model can be very different from the Programming Model**
  - E.g., von Neumann model implemented by an OoO processor
  - E.g., SPMD model implemented by a SIMD processor (a GPU)

# How Can You Exploit Parallelism Here?

```
for (i=0; i < N; i++)
```

*Scalar Sequential Code*

```
C[i] = A[i] + B[i];
```



Let's examine three programming options to exploit **instruction-level parallelism** present in this sequential code:

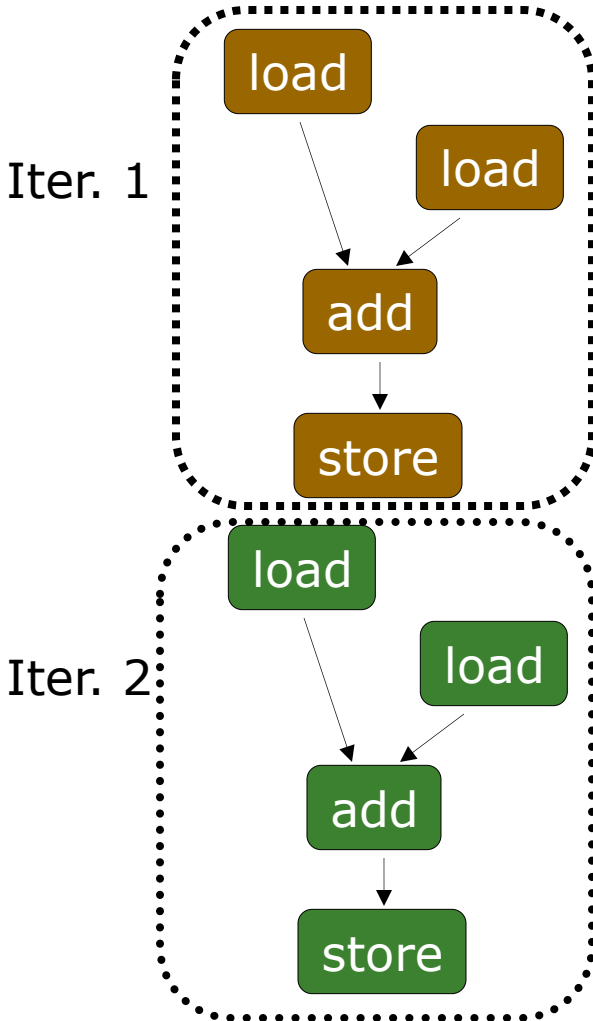
1. Sequential (SISD)
2. Data-Parallel (SIMD)
3. Multithreaded (MIMD/SPMD)



# Prog. Model 1: Sequential (SISD)

```
for (i=0; i < N; i++)  
    C[i] = A[i] + B[i];
```

## Scalar Sequential Code



- Can be executed on a:
  - Pipelined processor
  - Out-of-order execution processor
    - Independent instructions executed when ready
    - Different iterations are present in the instruction window and can execute in parallel in multiple functional units
    - In other words, the loop is dynamically unrolled by the hardware
- Superscalar or VLIW processor
  - Can fetch and execute multiple instructions per cycle

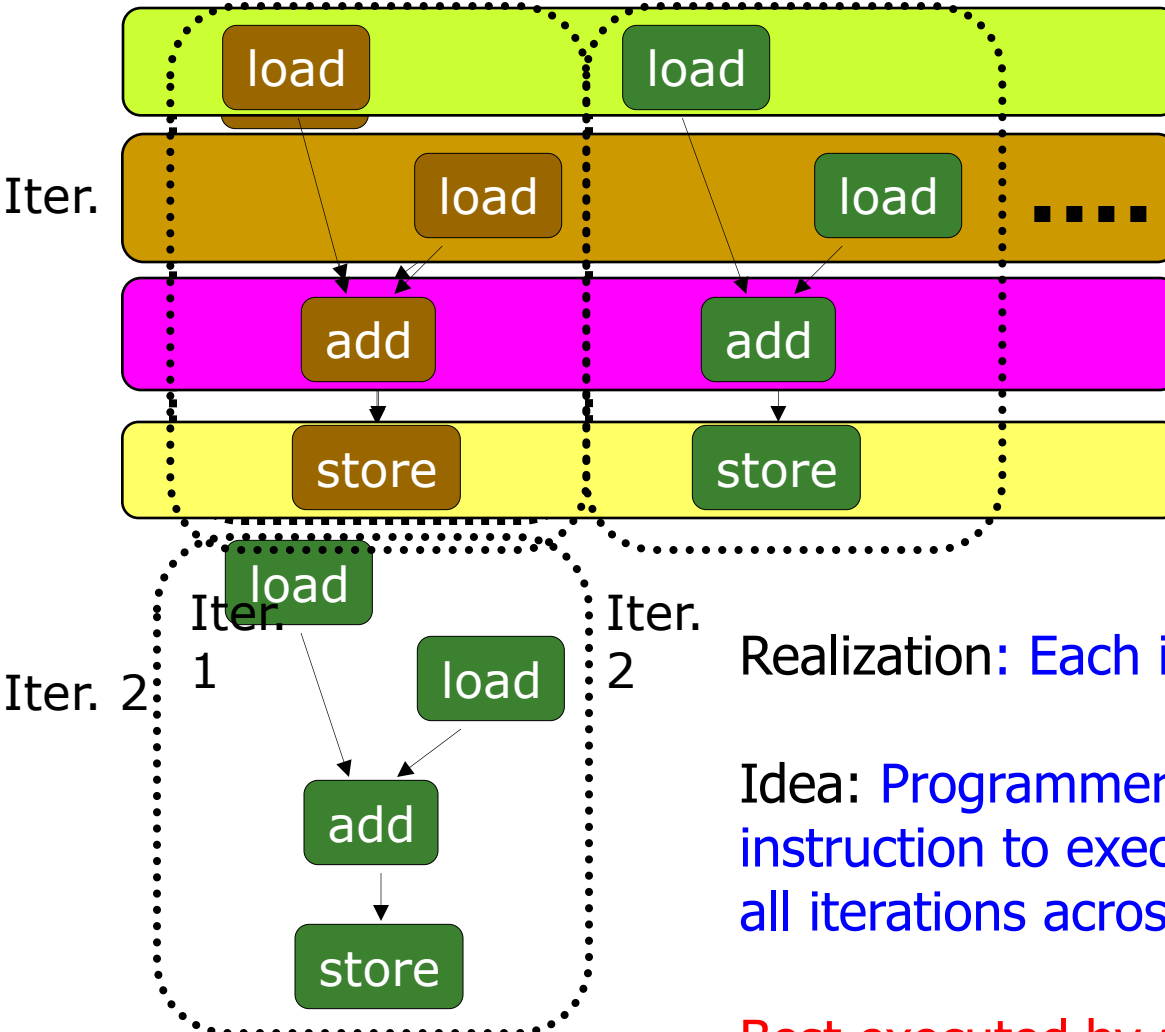
# Prog. Model 2: Data Parallel (SIMD)

```
for (i=0; i < N; i++)  
    C[i] = A[i] + B[i];
```

*Scalar Sequential Code*

*Vector Instruction*

*Vectorized Code*



VLD  $A \rightarrow V1$

VLD  $B \rightarrow V2$

VADD  $V1 + V2 \rightarrow V3$

VST  $V3 \rightarrow C$

Realization: Each iteration is independent

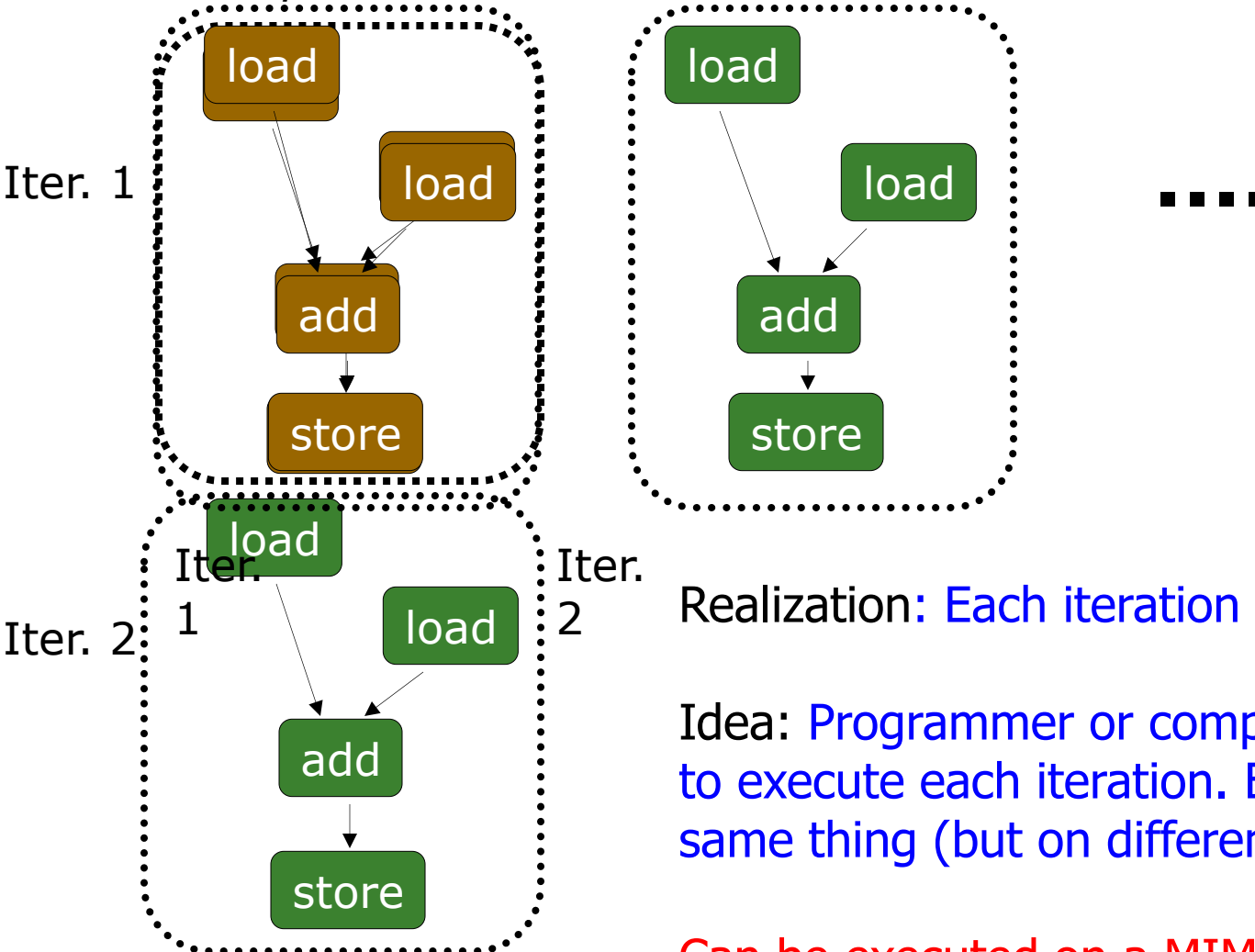
Idea: Programmer or compiler generates a SIMD instruction to execute the same instruction from all iterations across different data

Best executed by a SIMD processor (vector, array)

# Prog. Model 3: Multithreaded

```
for (i=0; i < N; i++)  
    C[i] = A[i] + B[i];
```

## Scalar Sequential Code



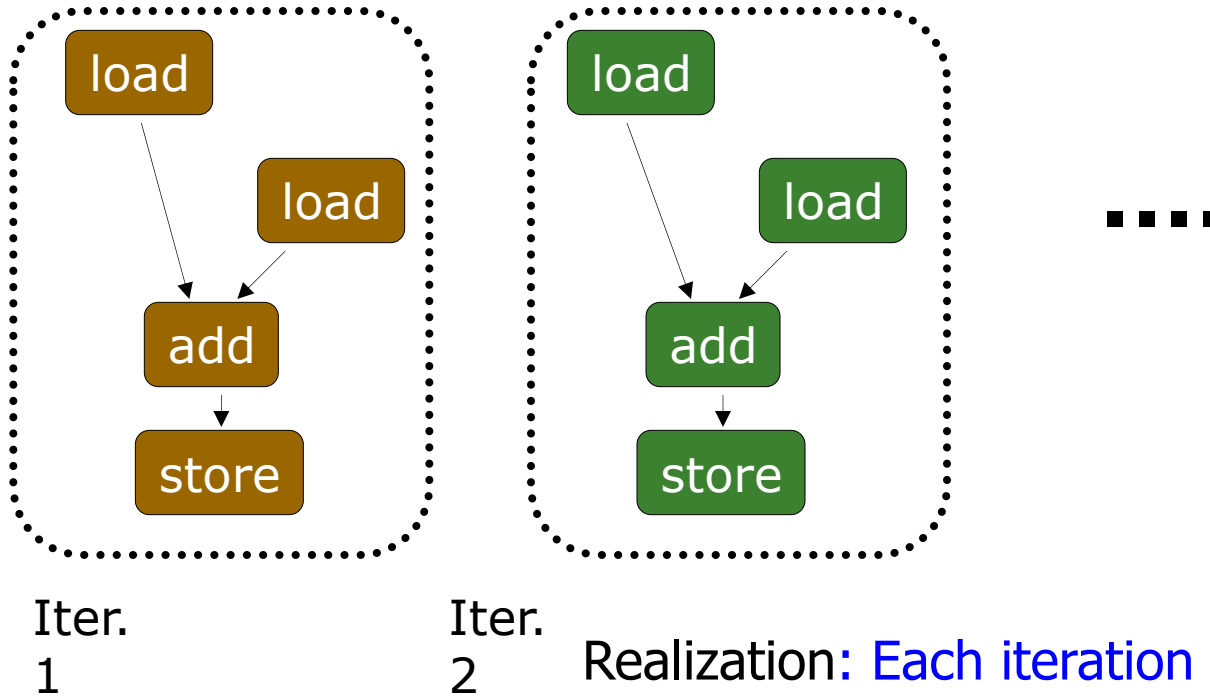
Realization: Each iteration is independent

Idea: Programmer or compiler generates a thread to execute each iteration. Each thread does the same thing (but on different data)

Can be executed on a MIMD machine

# Prog. Model 3: Multithreaded

```
for (i=0; i < N; i++)  
    C[i] = A[i] + B[i];
```



This particular model is also called:

**SPMD: Single Program Multiple Data**

Can be executed on a SIMT machine

**Single Instruction Multiple Thread**

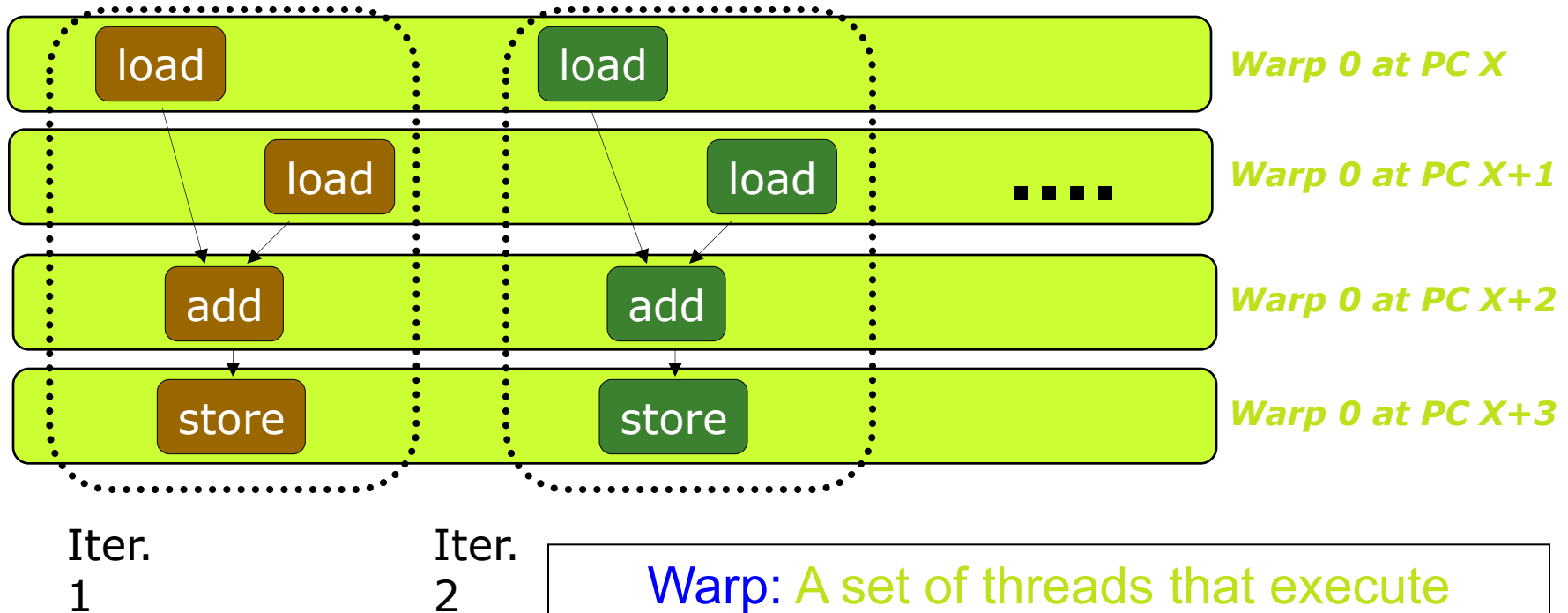
# A GPU is a SIMD (SIMT) Machine

---

- Except it is **not** programmed using SIMD instructions
- It is **programmed using threads** (SPMD programming model)
  - Each thread executes the same code but operates on a different piece of data
  - Each thread has its own context (i.e., can be treated/restarted/executed independently)
- A set of threads executing the same instruction are dynamically grouped into a **warp (wavefront)** by the hardware
  - A warp is essentially a **SIMD operation formed by hardware!**

# SPMD on SIMT Machine

```
for (i=0; i < N; i++)  
    C[i] = A[i] + B[i];
```



**Warp:** A set of threads that execute the same instruction (i.e., at the same PC)

This particular model is also called:

**SPMD: Single Program Multiple Data**

A GPU executes it using the SIMT model:  
**Single Instruction Multiple Thread**

# Graphics Processing Units

SIMD not Exposed to Programmer (SIMT)

# SIMD vs. SIMT Execution Model

---

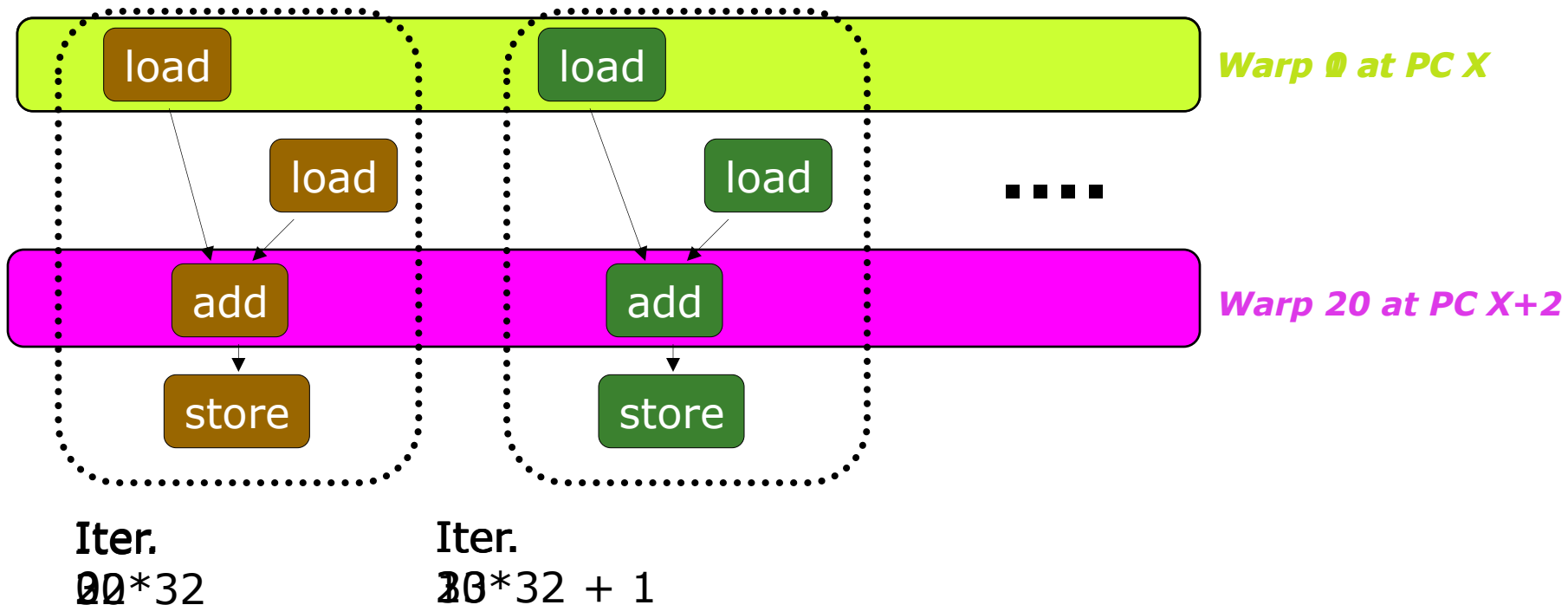
- SIMD: A single **sequential instruction stream** of **SIMD instructions** → each instruction specifies multiple data inputs
  - [VLD, VLD, VADD, VST], VLEN
- SIMT: **Multiple instruction streams** of **scalar instructions** → threads grouped dynamically into warps
  - [LD, LD, ADD, ST], NumThreads
- Two Major SIMT Advantages:
  - **Can treat each thread separately** → i.e., can execute each thread independently (on any type of scalar pipeline) → MIMD processing
  - **Can group threads into warps flexibly** → i.e., can group threads that are supposed to *truly* execute the same instruction → dynamically obtain and maximize benefits of SIMD processing



# Fine-Grained Multithreading of Warps

```
for (i=0; i < N; i++)  
    C[i] = A[i] + B[i];
```

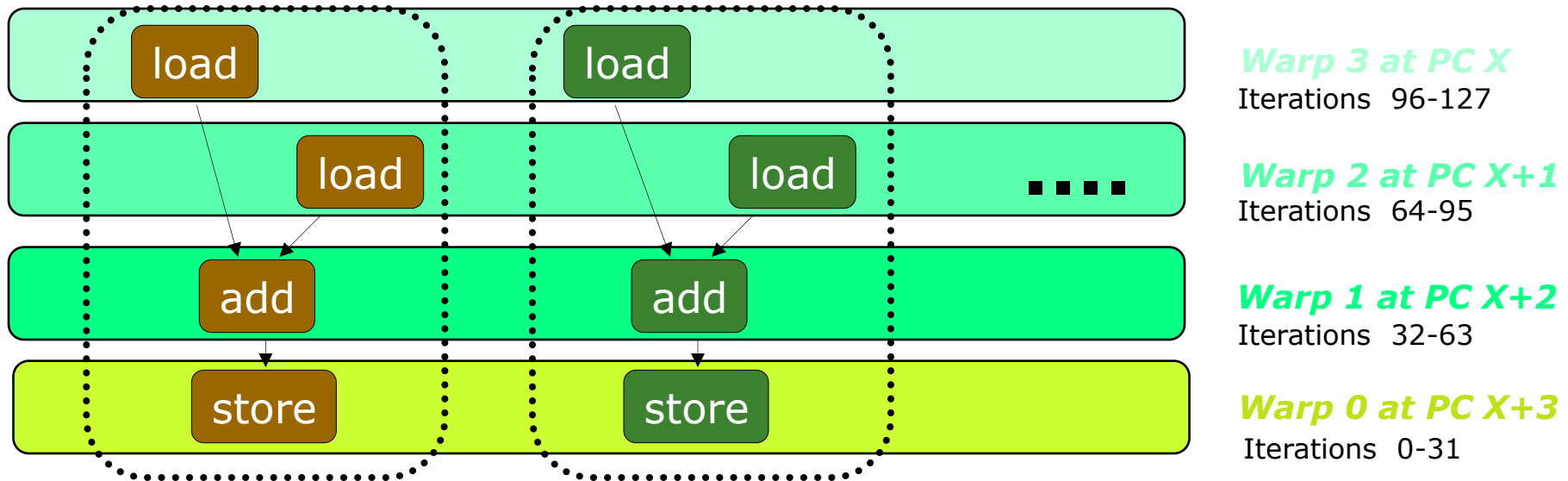
- Assume a warp consists of 32 threads
- If you have 32K iterations, and 1 iteration/thread  $\rightarrow$  1K warps
- Warps can be interleaved on the same pipeline  $\rightarrow$  Fine grained multithreading of warps



# Fine-Grained Multithreading of Warps

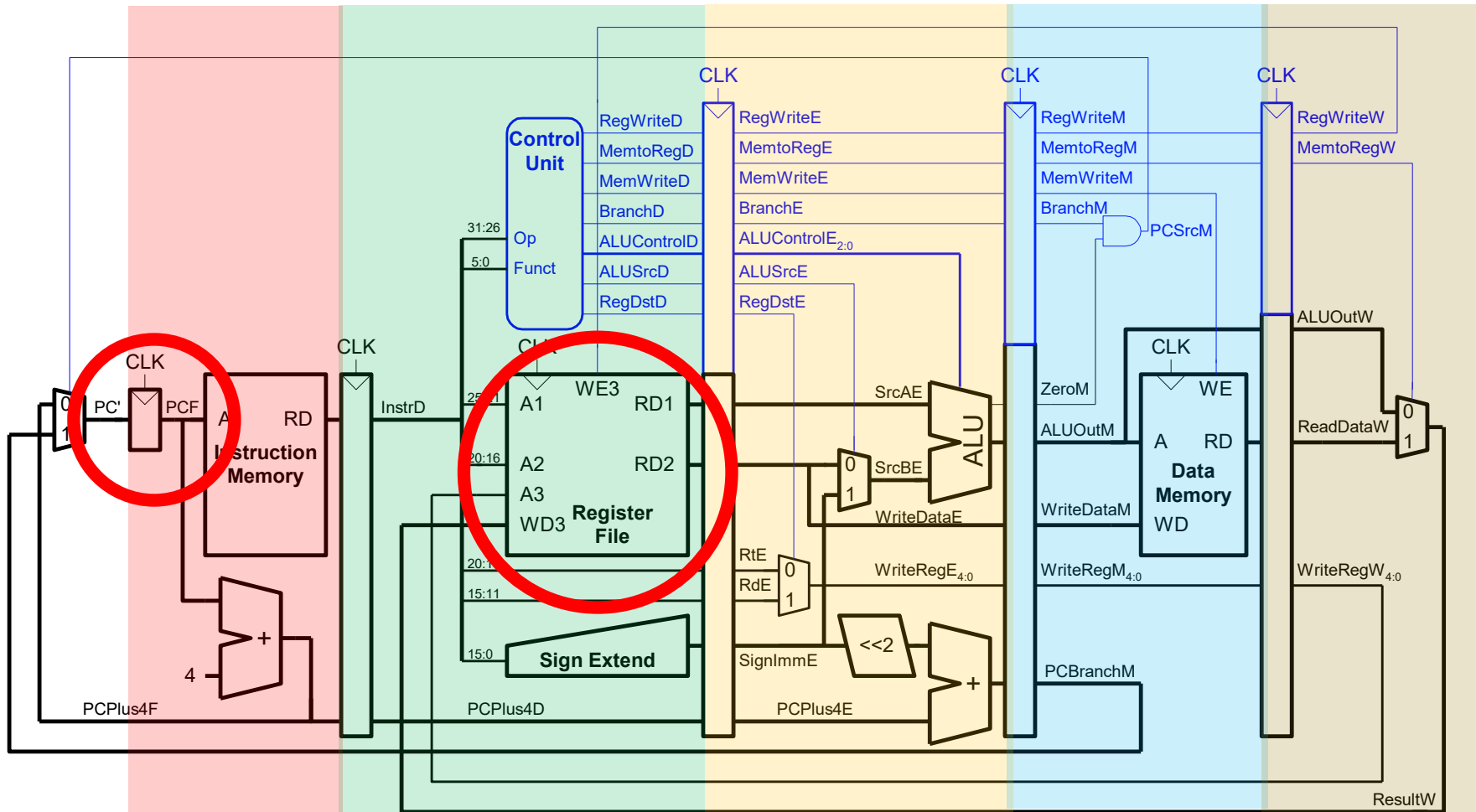
```
for (i=0; i < N; i++)  
    C[i] = A[i] + B[i];
```

- Assume a warp consists of 32 threads
- If you have 32K iterations, and 1 iteration/thread → 1K warps
- Warps can be interleaved on the same pipeline → Fine grained multithreading of warps



All threads in a warp are independent of each other  
→ They be executed seamlessly in a fine-grained multithreaded pipeline

# Recall: Fine-Grained Multithreading: Basic Idea

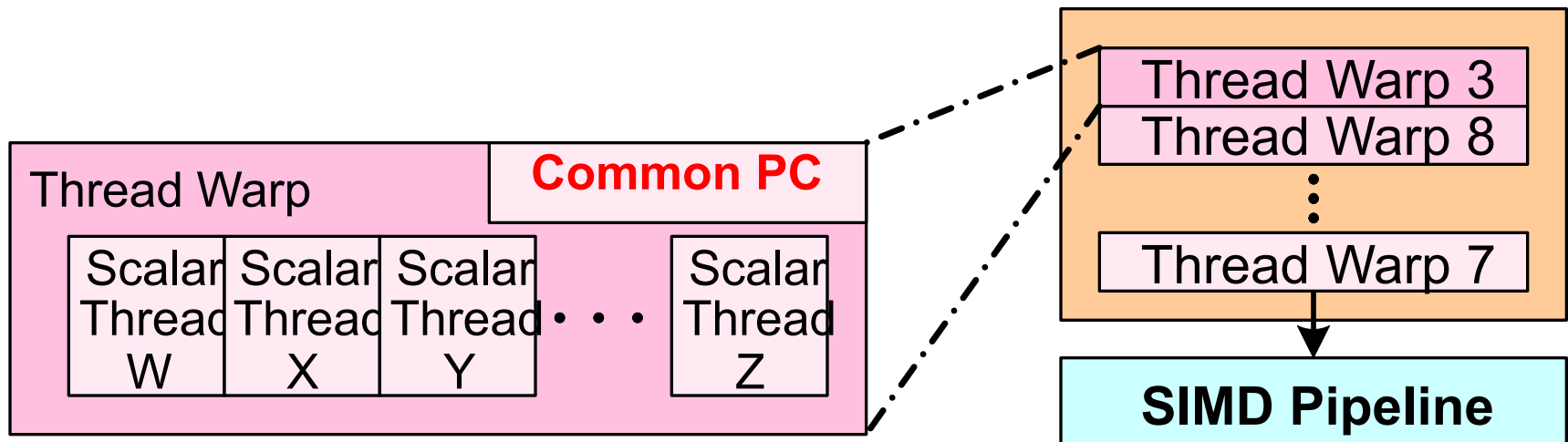


**Each pipeline stage has an instruction from a different, completely-independent thread**

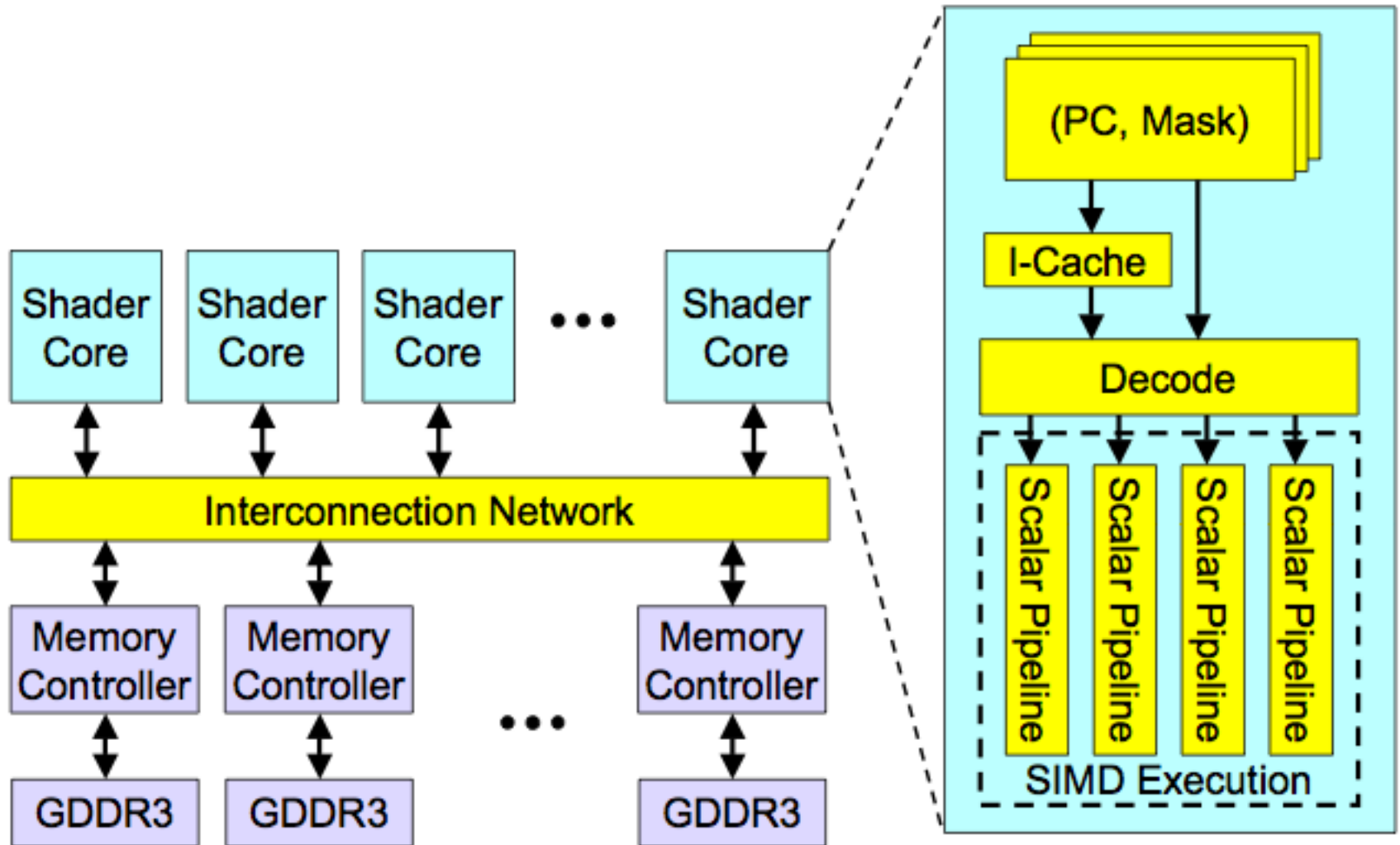
**We need a PC and a register file for each thread + muxes and control**

# Warps and Warp-Level FGMT

- Warp: A set of threads that execute the same instruction (on different data elements) → SIMT (Nvidia-speak)
- All threads run the same code
- Warp: The threads that run lengthwise in a woven fabric ...

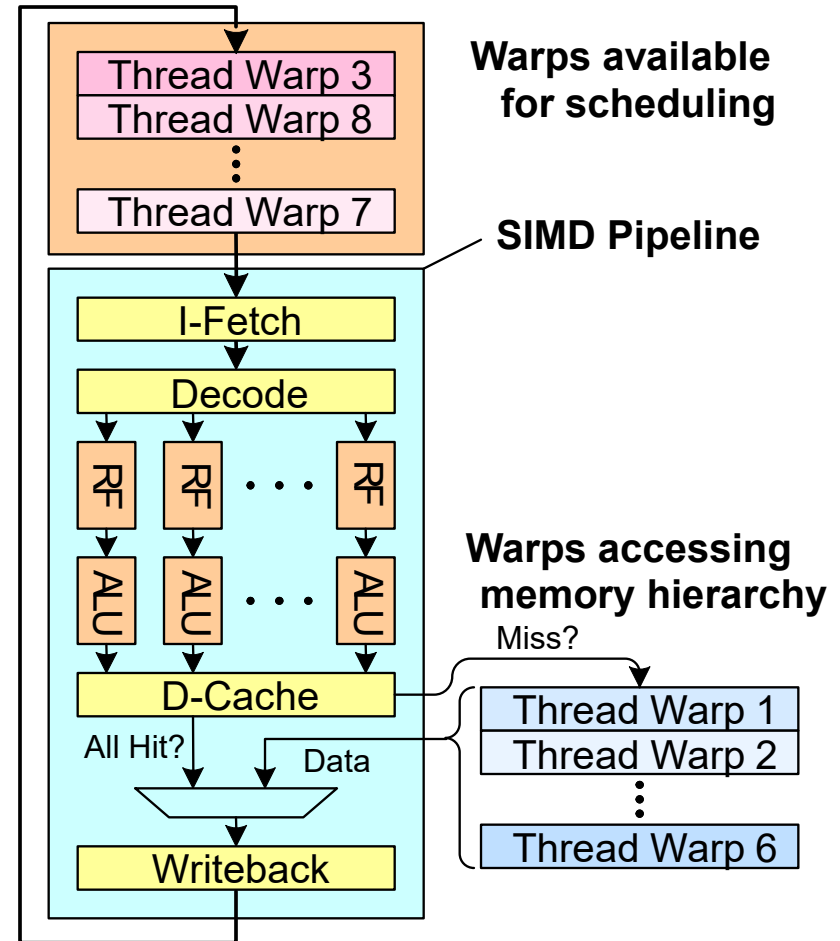


# High-Level View of a GPU



# 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 simple pipeline & long latency tolerance
  - ❑ Millions of threads operating on the same large image/video

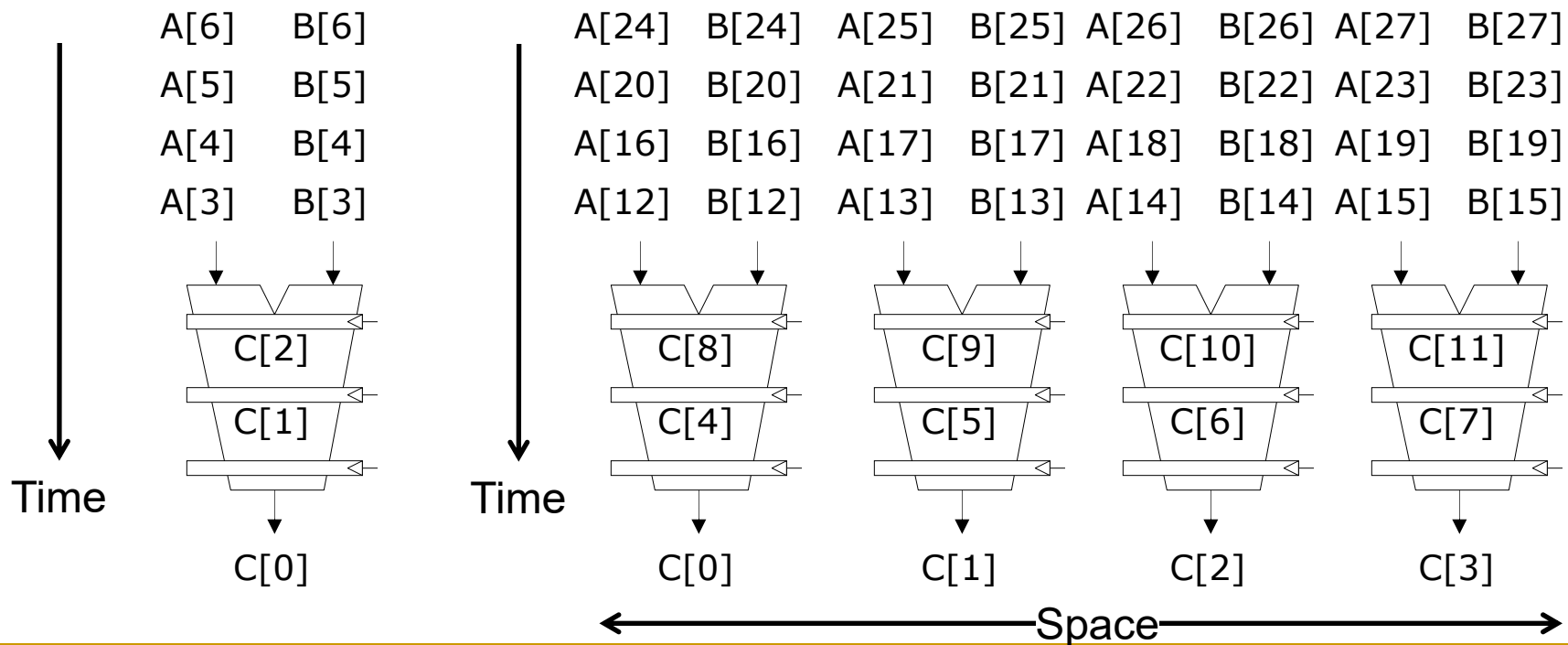


# Recall: Vector Instruction Execution

**VADD A,B → C**

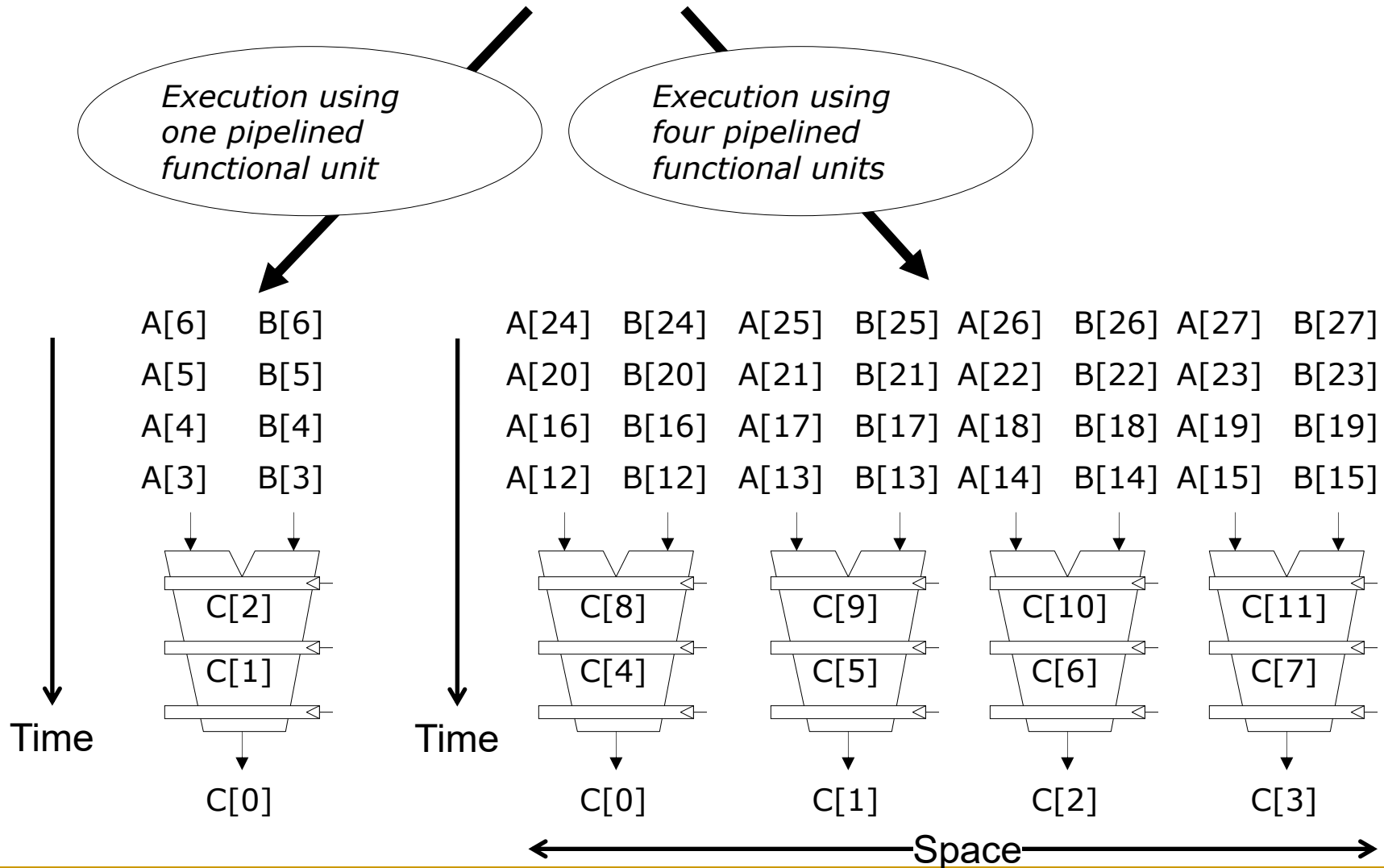
*Execution using  
one pipelined  
functional unit*

*Execution using  
four pipelined  
functional units*



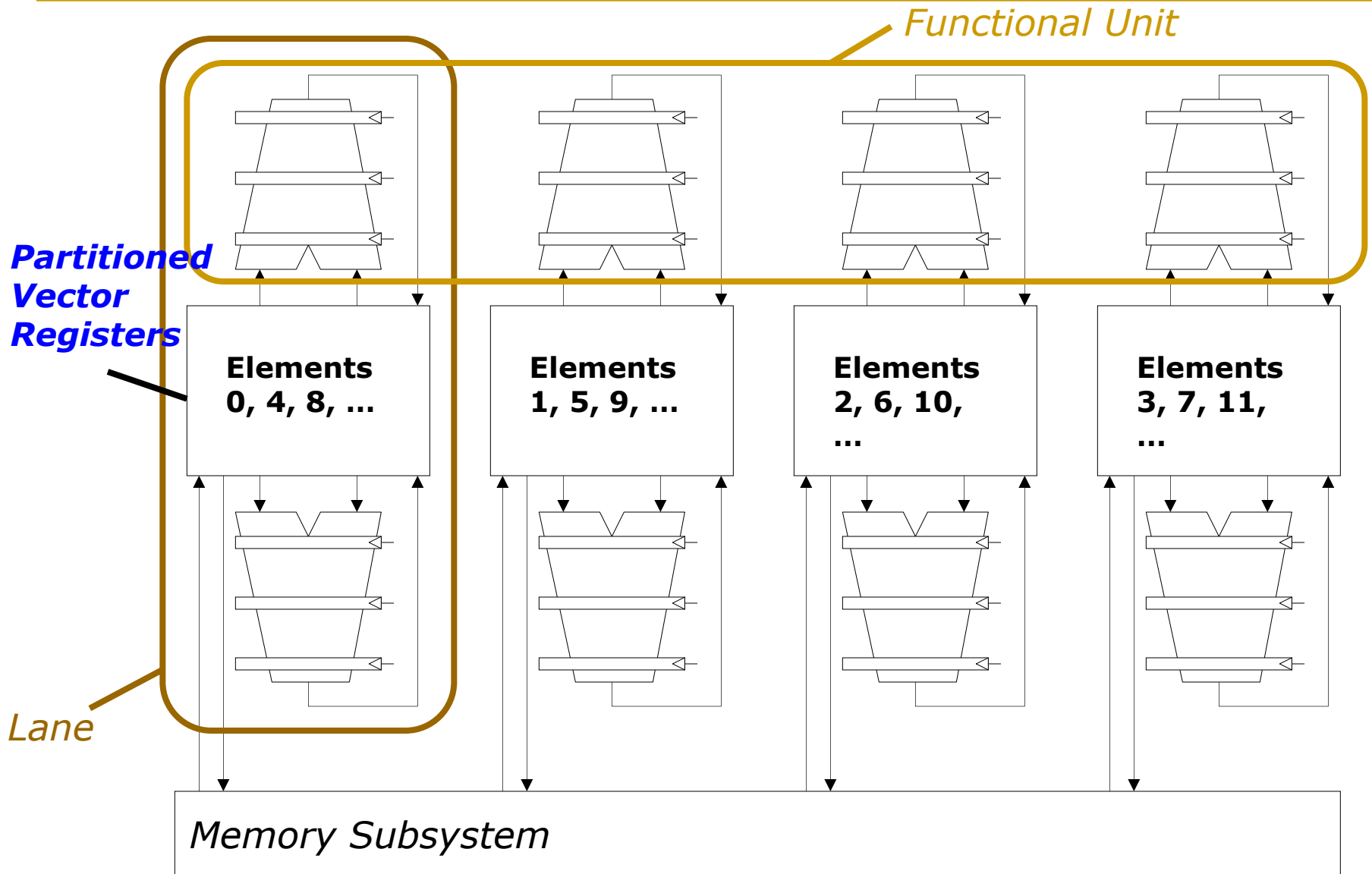
# Warp Execution (Recall the Previous Slide)

32-thread warp executing **ADD A[tid],B[tid] → C[tid]**

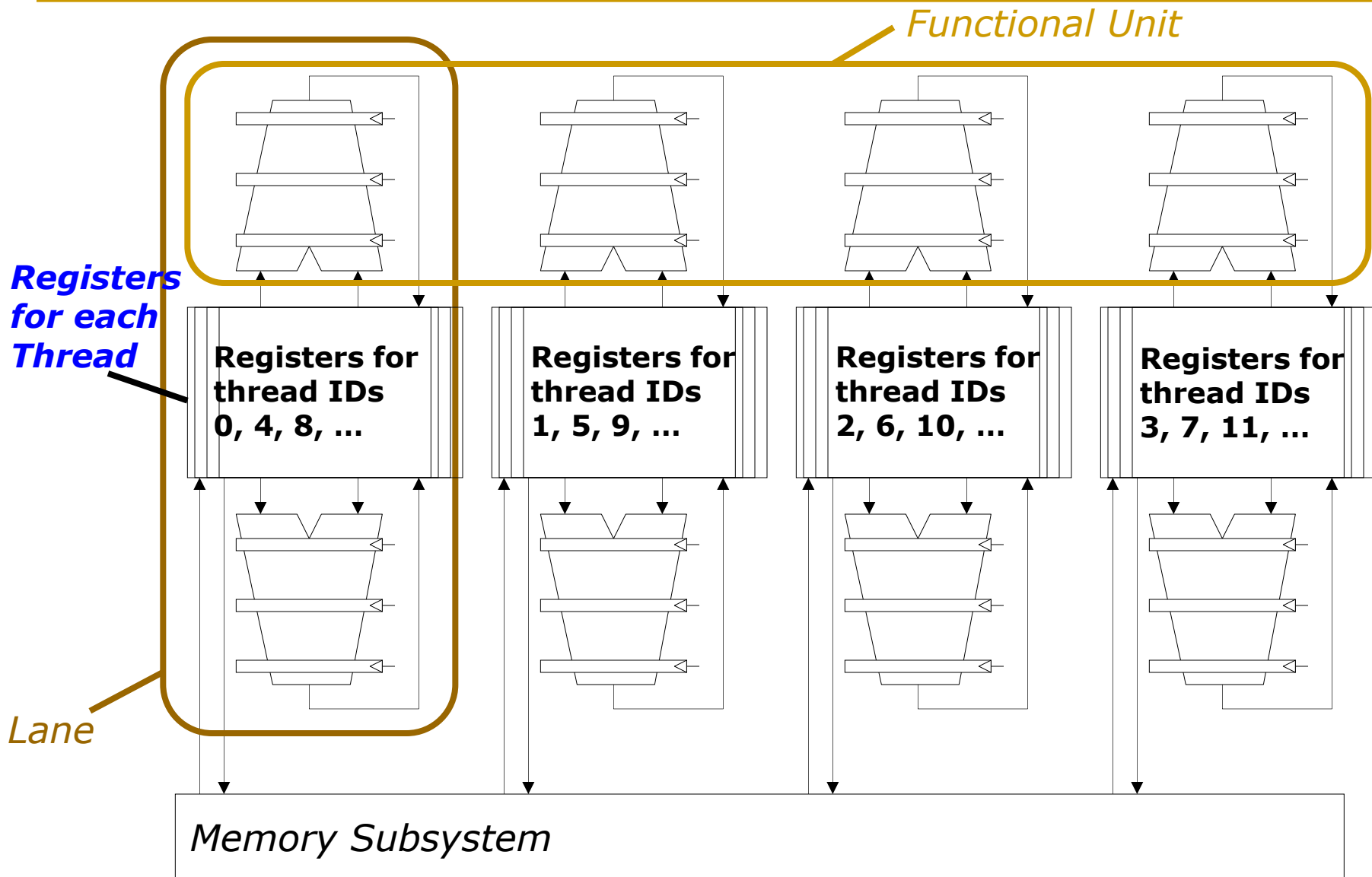




# Recall: Vector Unit Structure



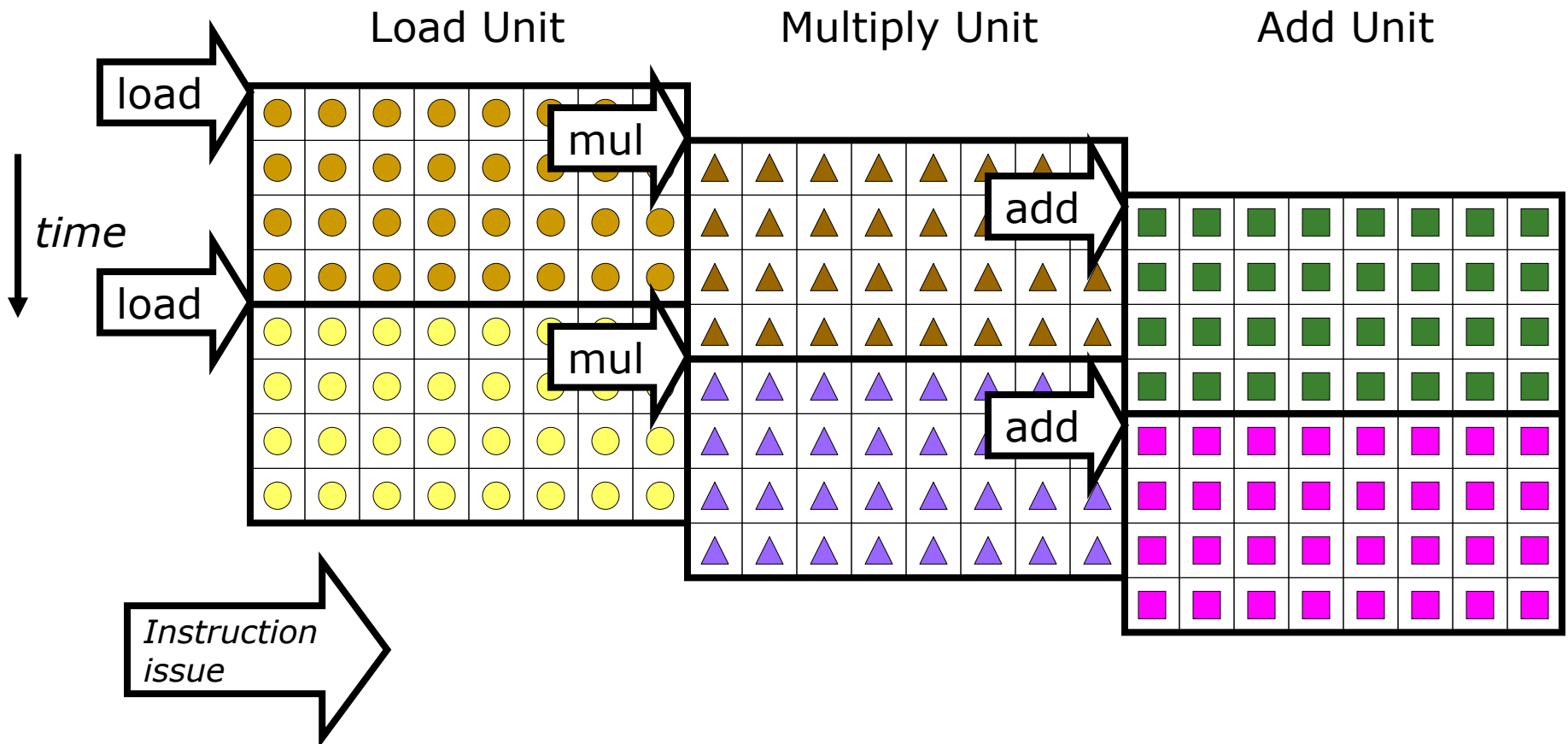
# GPU SIMD Execution Unit Structure



# Recall: Vector Instruction Level Parallelism

Can overlap execution of multiple vector instructions

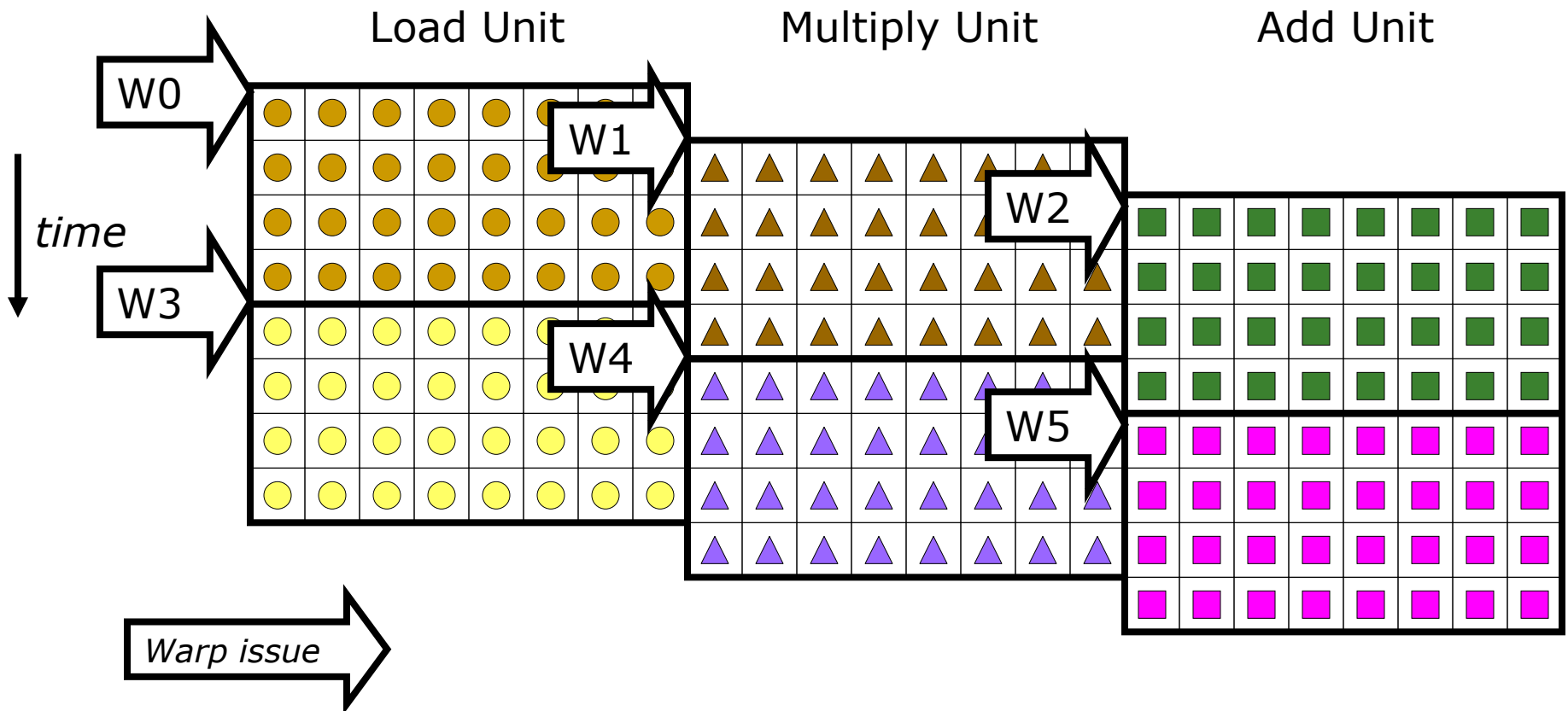
- Example machine has 32 elements per vector register and 8 lanes
- Example with 24 operations/cycle (steady state) while issuing 1 vector instruction/cycle



# Warp Instruction Level Parallelism

Can overlap execution of multiple instructions

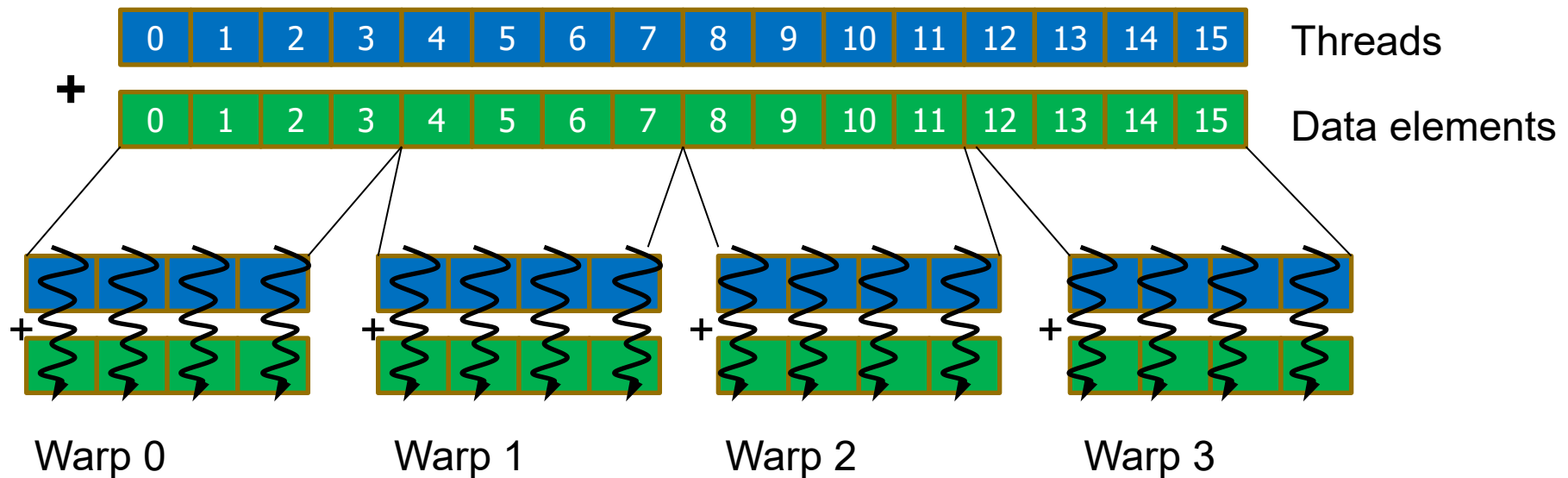
- Example machine has 32 threads per warp and 8 lanes
- Completes 24 operations/cycle (steady state) while issuing 1 warp/cycle



# SIMT Memory Access (Loads and Stores)

- Same instruction in different threads uses **thread id** to index and access different data elements

Let's assume  $N=16$ , 4 threads per warp  $\rightarrow$  4 warps



**For maximum performance, memory should provide enough bandwidth**  
(i.e., elements per cycle throughput to match computation unit throughput)

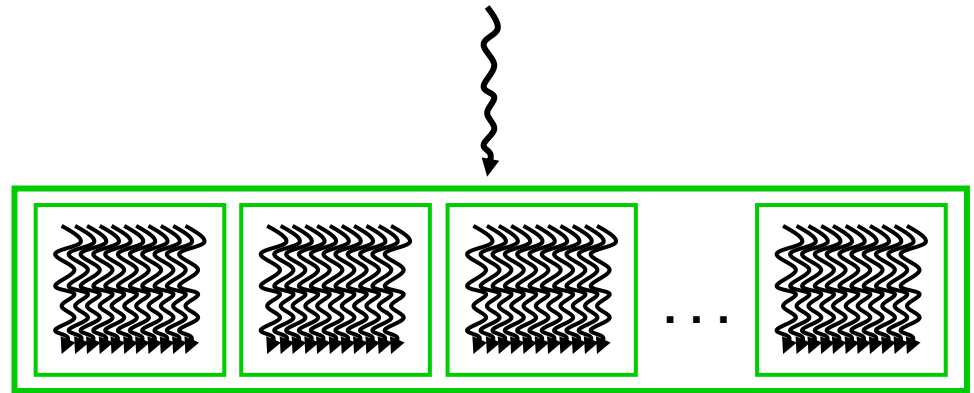
# Warps *not* Exposed to GPU Programmers

- CPU threads and GPU kernels
  - ▣ Sequential or modestly parallel sections on CPU
  - ▣ Massively parallel sections on GPU: Blocks of threads

Serial Code (host)

Parallel Kernel (device)

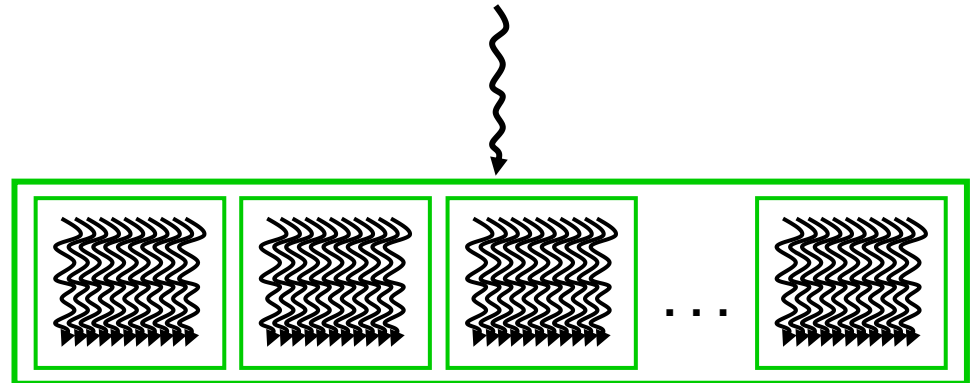
```
KernelA<<<nBlk, nThr>>>(args);
```



Serial Code (host)

Parallel Kernel (device)

```
KernelB<<<nBlk, nThr>>>(args);
```



# Sample GPU SIMT Code (Simplified)

---

CPU code

```
for (ii = 0; ii < 100000; ++ii) {  
    C[ii] = A[ii] + B[ii];  
}
```



CUDA code

```
// there are 100000 threads  
__global__ void KernelFunction(...) {  
    int tid = blockDim.x * blockIdx.x + threadIdx.x;  
    int varA = aa[tid];  
    int varB = bb[tid];  
    C[tid] = varA + varB;  
}
```

# Sample GPU Program (Less Simplified)

## CPU Program

```
void add matrix
( float *a, float* b, float *c, int N) {
    int index;
    for (int i = 0; i < N; ++i)
        for (int j = 0; j < N; ++j) {
            index = i + j*N;
            c[index] = a[index] + b[index];
        }
}

int main () {

    add matrix (a, b, c, N);
}
```

## GPU Program

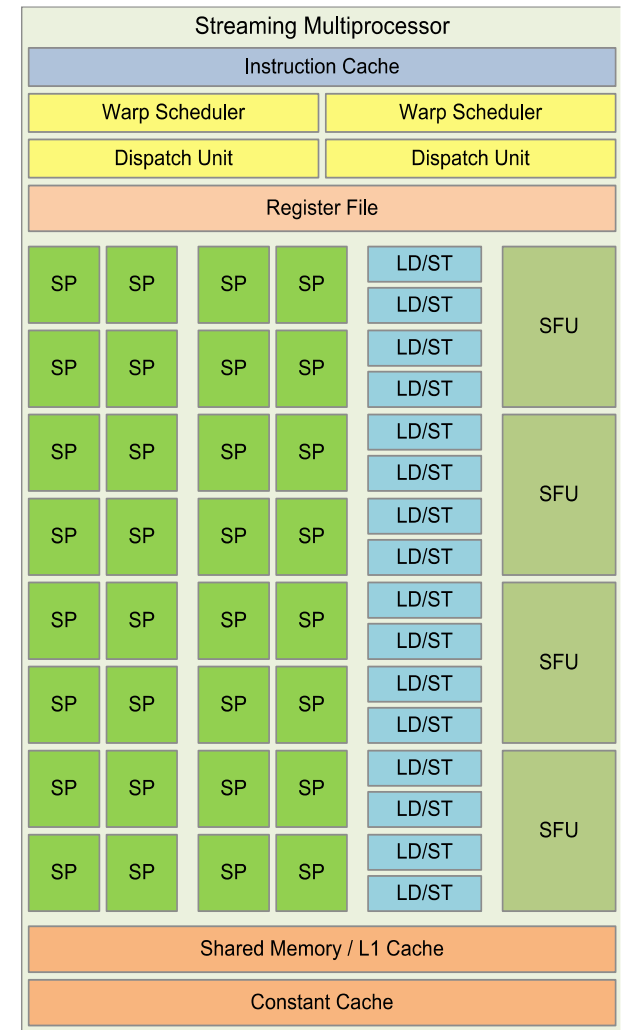
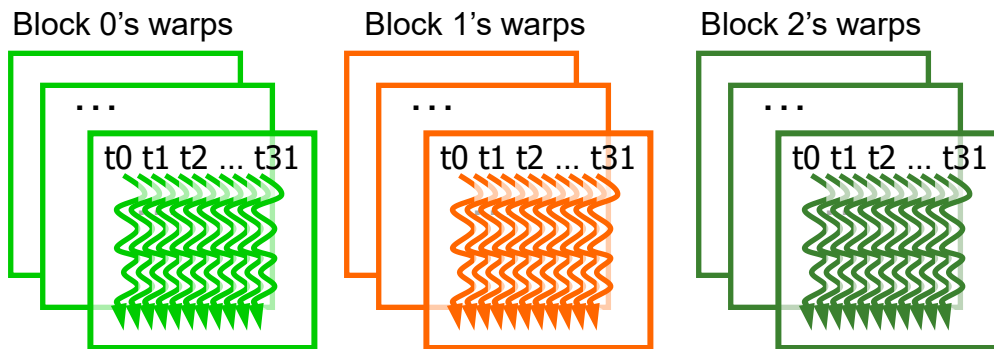
```
__global__ add_matrix
( float *a, float *b, float *c, int N) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;
    int index = i + j*N;
    if (i < N && j < N)
        c[index] = a[index]+b[index];
}

int main() {
    dim3 dimBlock( blocksize, blocksize) ;
    dim3 dimGrid (N/dimBlock.x, N/dimBlock.y);
    add_matrix<<<dimGrid, dimBlock>>>( a, b, c, N);
}
```



# From Blocks to Warps

- GPU core: A SIMD pipeline
  - ❑ Streaming Processor (SP)
  - ❑ Many such SIMD Processors
    - Streaming Multiprocessor (SM)
- Blocks are divided into **warps**
  - ❑ SIMD/SIMT unit (32 threads)



NVIDIA Fermi architecture

# Warp-based SIMD vs. Traditional SIMD

---

- **Traditional SIMD** contains a **single thread**
  - ❑ **Sequential instruction execution**; lock-step operations in a SIMD instruction
  - ❑ **Programming model is SIMD** (no extra threads) → SW needs to know vector length
  - ❑ ISA contains **vector/SIMD instructions**
- **Warp-based SIMD** consists of **multiple scalar threads** executing in a SIMD manner (i.e., same instruction executed by all threads)
  - ❑ Does not have to be lock step
  - ❑ **Each thread can be treated individually** (i.e., placed in a different warp) → **programming model not SIMD**
    - SW does **not need to know vector length**
    - Enables multithreading and flexible dynamic grouping of threads
  - ❑ **ISA is scalar** → SIMD operations can be formed dynamically
  - ❑ Essentially, it is **SPMD programming model implemented on SIMD hardware**

# SPMD

---

- Single procedure/program, multiple data
  - This is a programming model rather than computer organization
- Each processing element executes the same procedure, except on different data elements
  - Procedures can synchronize at certain points in program, e.g. barriers
- Multiple instruction streams execute the same program
  - Each program/procedure 1) works on different data, 2) can execute a different control-flow path, at run-time
  - Many scientific applications are programmed this way and run on MIMD hardware (multiprocessors)
  - Modern GPUs programmed in a similar way on a SIMD hardware

# SIMD vs. SIMT Execution Model

---

- SIMD: A single **sequential instruction stream** of **SIMD instructions** → each instruction specifies multiple data inputs
  - [VLD, VLD, VADD, VST], VLEN
- SIMT: **Multiple instruction streams** of **scalar instructions** → threads grouped dynamically into warps
  - [LD, LD, ADD, ST], NumThreads
- Two Major SIMT Advantages:
  - **Can treat each thread separately** → i.e., can execute each thread independently on any type of scalar pipeline → MIMD processing
  - **Can group threads into warps flexibly** → i.e., can group threads that are supposed to *truly* execute the same instruction → dynamically obtain and maximize benefits of SIMD processing

# TSPs (Tensor Streaming Processors)

# The Tensor Streaming Processor (TSP)

---

- **Definition:** A novel processing architecture designed specifically for deep learning and linear algebra operations.
- **Philosophy:** Moves complexity from hardware to software. The compiler plans the exact execution of every instruction before the chip runs it.
- **Primary Goal:** To solve the "latency vs. throughput" trade-off common in GPUs.

# Key Architectural Diffs (TSP vs. GPU)

---

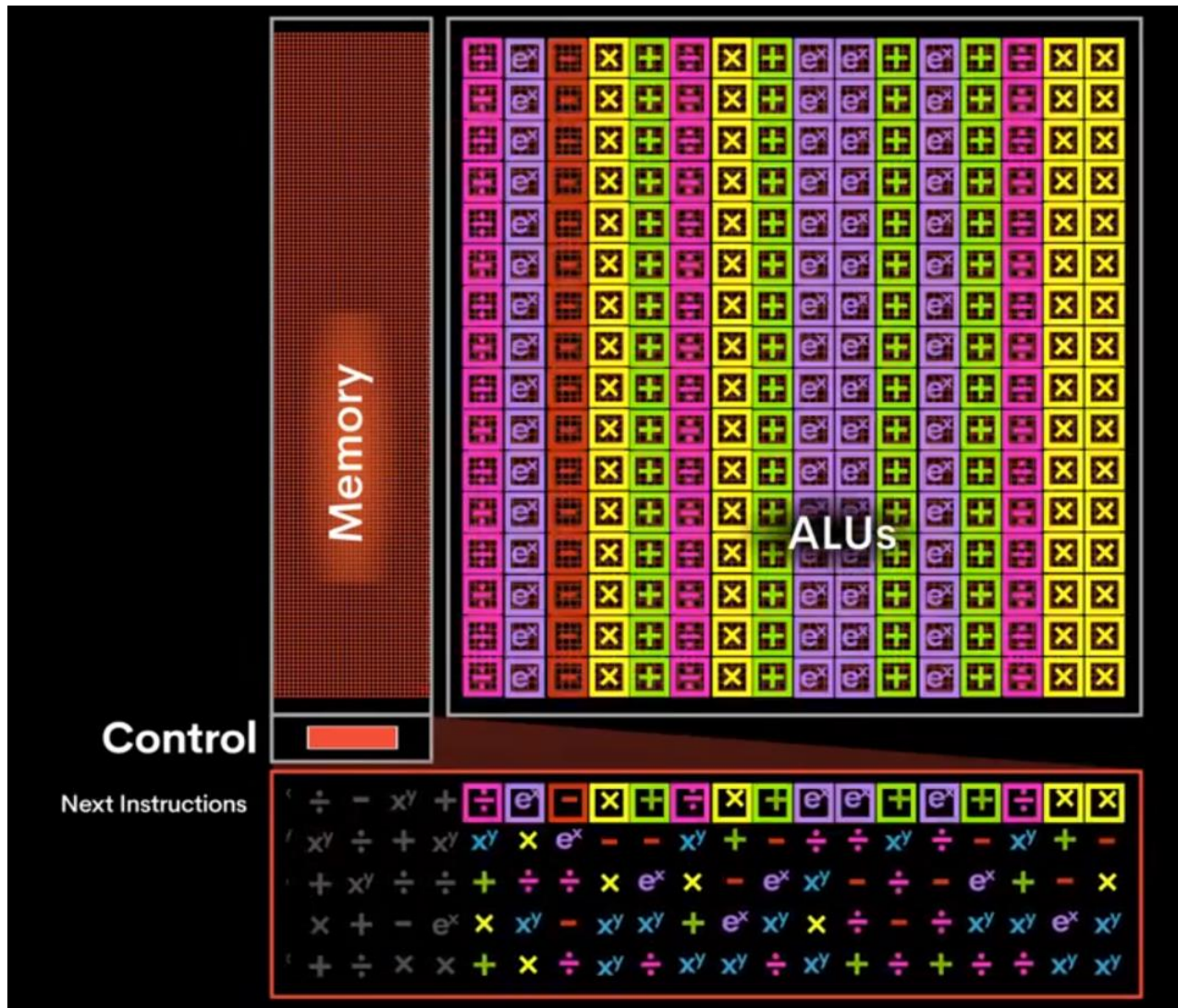
Feature	GPU (Graphics Processing Unit)	TSP (Tensor Streaming Processor)
Control	Hardware-managed (schedulers, arbiters)	<b>Software-defined</b> (Compiler-managed)
Data Flow	Complex memory hierarchy (L1/L2 caches)	<b>Stream-based</b> (Direct register access)
Execution	Non-deterministic (Reactive)	<b>Deterministic</b> (Predictable)
Cores	Thousands of small cores	Single, massive monolithic core

# The "Assembly Line" Concept

---

- **Think of a TSP like a factory assembly line.** Data enters the chip and moves in a perfectly timed, rhythmic flow across functional units (Vector, Matrix, Switch). There is no traffic, no waiting, and no "guessing" (branch prediction).
- **Determinism:** The compiler knows *exactly* when data will arrive at a specific unit. Explain that on a GPU, performance varies based on cache hits and thread scheduling. On a TSP, if the compiler says a task takes 500 microseconds, it takes *exactly* 500 microseconds every single time.
- **No Caches:** Removes cache misses and memory stalls.
- **Batch Size 1 Performance:** Delivers high throughput even without batching requests, resulting in ultra-low latency for Real-Time AI.





# Thank you.

---