

High Performance Computing for Science and Engineering II

29.4.2019 - **Lecture 9: CUDA Optimization**

Lecturer: Dr. Sergio Martin

CSElab

Computational Science & Engineering Laboratory

ETH zürich

General Info

Semester Exam

- **When:** May 27th – 10.00am – 1.00pm
- **Where:** HG E7
- **What:** 3 Hours - Handwritten Only
- Can contain code-related questions.
- We'll let you know what aids you can bring.

Next Week's Exercise Session

- We will have another interactive session (CUDA)
- Bring your laptop.

Today's Class

A brief history of memory hierarchies.

- From first cache structures modern CPUs

CUDA: Optimizing Memory Access

- Shared memory, register memory. Example.

Kernel Decomposition

- Vector reduction pattern.

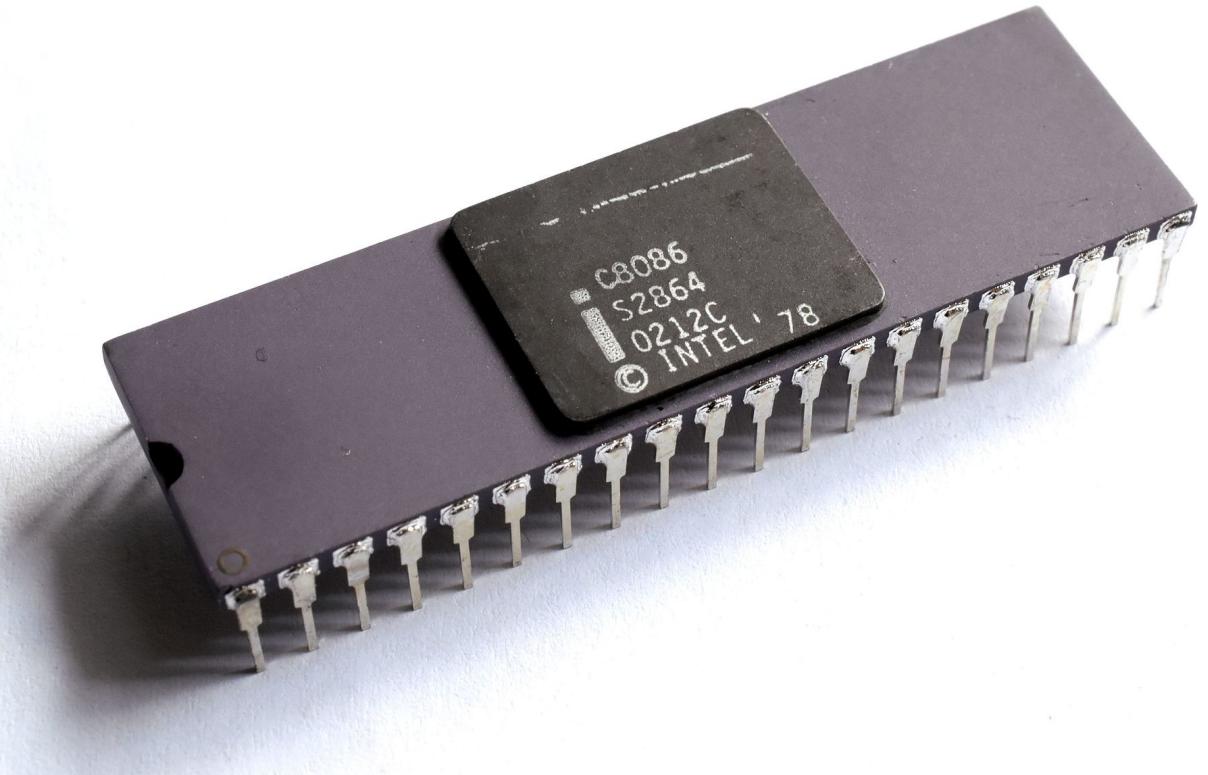
Homework 5

- Goals and Grading.

Hardware Hierarchies: A brief History

The Intel 8086 Processor

1978 - Intel Releases 8086, the first 16-bit processor of the x86 architecture.

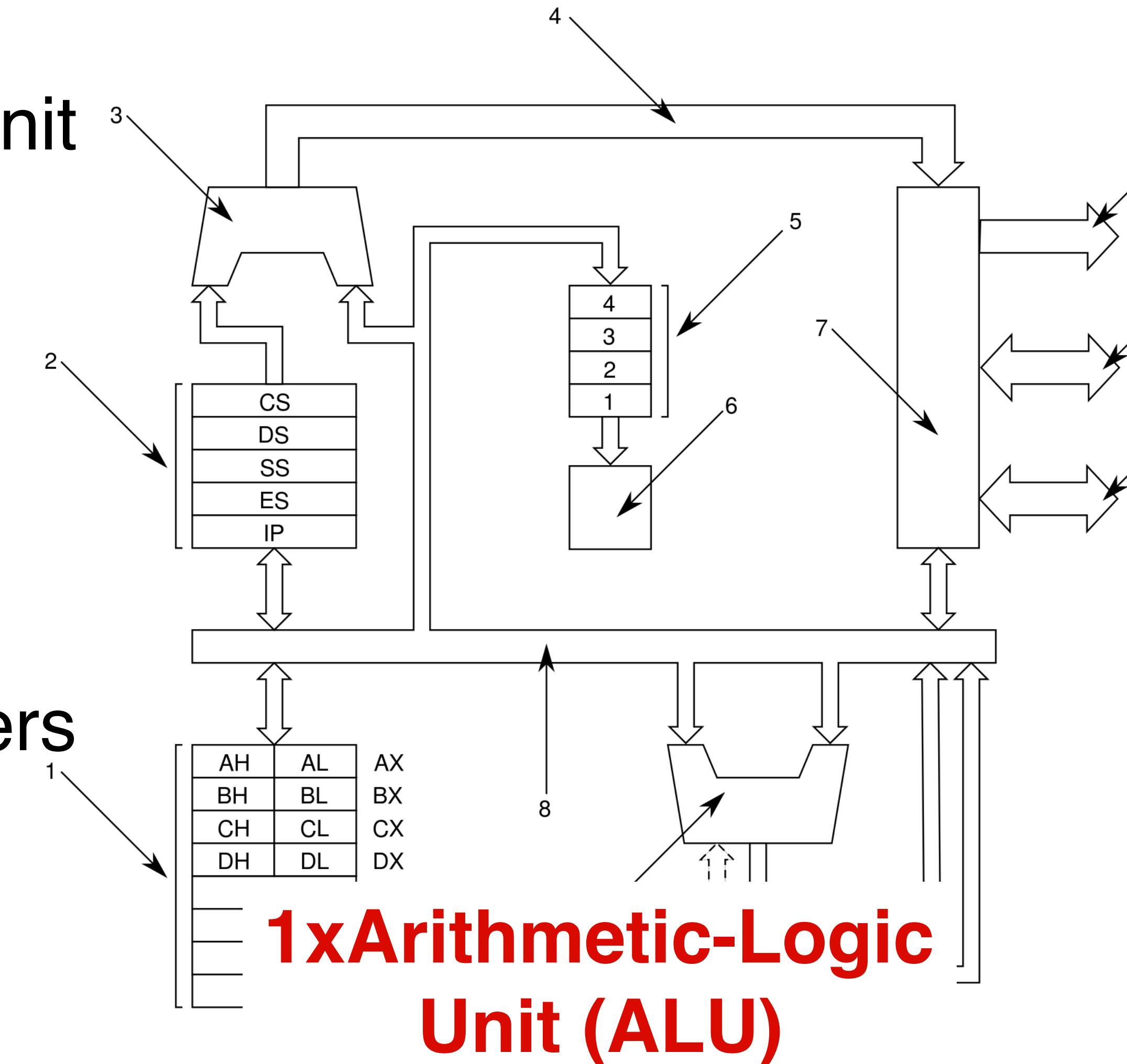


29k Transistors - 5Mhz

16-bit = 64k

Max Memory Address

Address Calculator Unit
Segment Selectors (16-bit)
General-Purpose Registers (16-bit)



1x Arithmetic-Logic Unit (ALU)

(32-bit) Extended Registers

1986 - Intel Releases 80386, its first 32-bit Processor.



275k Transistors - 20Mhz

32-bit = 4gb

Max Memory Address

Extended
Registers

Applications Operate
on Larger Data Sets

Additional Pressure
on RAM

Intel 80386 registers			
3 ₁	...	1 ₅	...
(bit position)			
Main registers (8/16/32 bits)			
EAX	AX	AL	Accumulator register
EBX	BX	BL	Base register
ECX	CX	CL	Count register
EDX	DX	DL	Data register
Index registers (16/32 bits)			
ESI	SI		Source Index
EDI	DI		Destination Index
EBP	BP		Base Pointer
ESP	SP		Stack Pointer
Program counter (16/32 bits)			
EIP	IP		Instruction Pointer
Segment selectors (16 bits)			
	CS		Code Segment
	DS		Data Segment
	ES		Extra Segment
	FS		F Segment
	GS		G Segment
	SS		Stack Segment
Status register			
1 ₇	1 ₆	1 ₅	1 ₄
1 ₃	1 ₂	1 ₁	1 ₀
0 ₉	0 ₈	0 ₇	0 ₆
0 ₅	0 ₄	0 ₃	0 ₂
0 ₁	0 ₀		
(bit position)			
V	R	0	N
I	O	P	L
D	I	S	Z
T	S	A	Z
S	Z	P	1
A	0	P	C
EFlags			

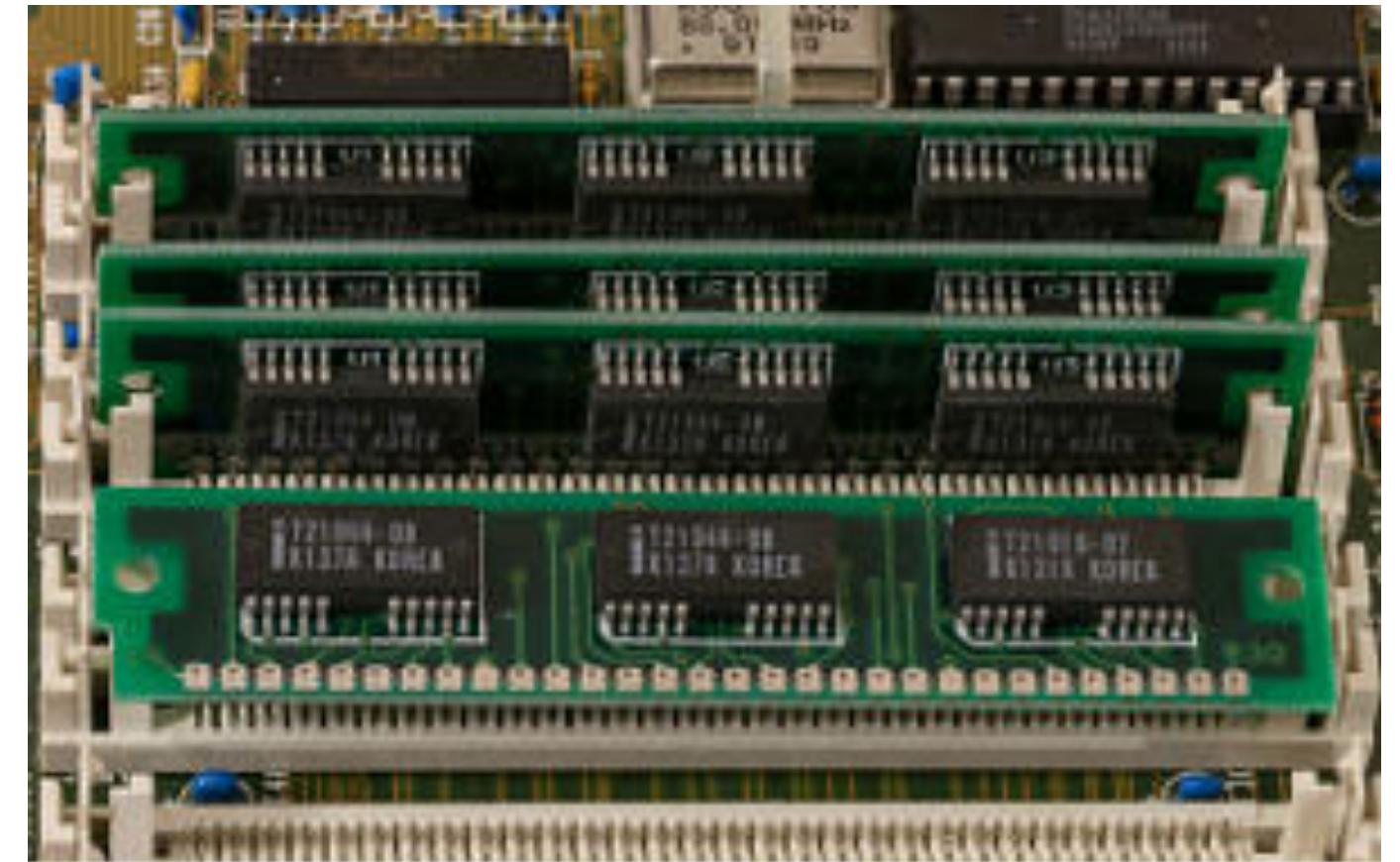
CPU/RAM Latency



275k Transistors - 20Mhz

GP Register Capacity: 16 Bytes

Register Memory Latency: 2-5ns



ComputerHope.com

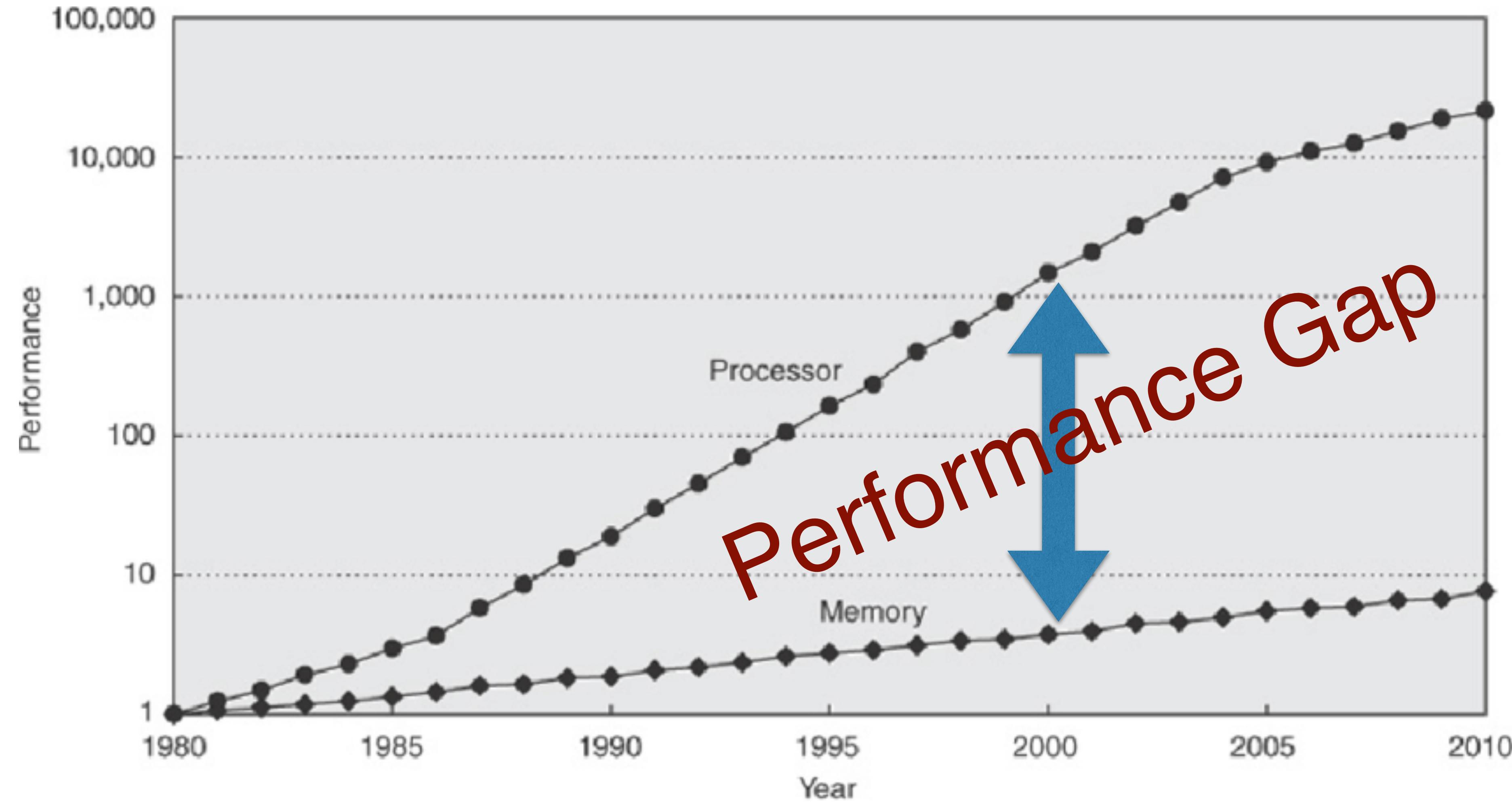
Asynchronous DRAM

Capacity: 1-64 Mbytes

Latency: ~120ns

Memory Performance Gap

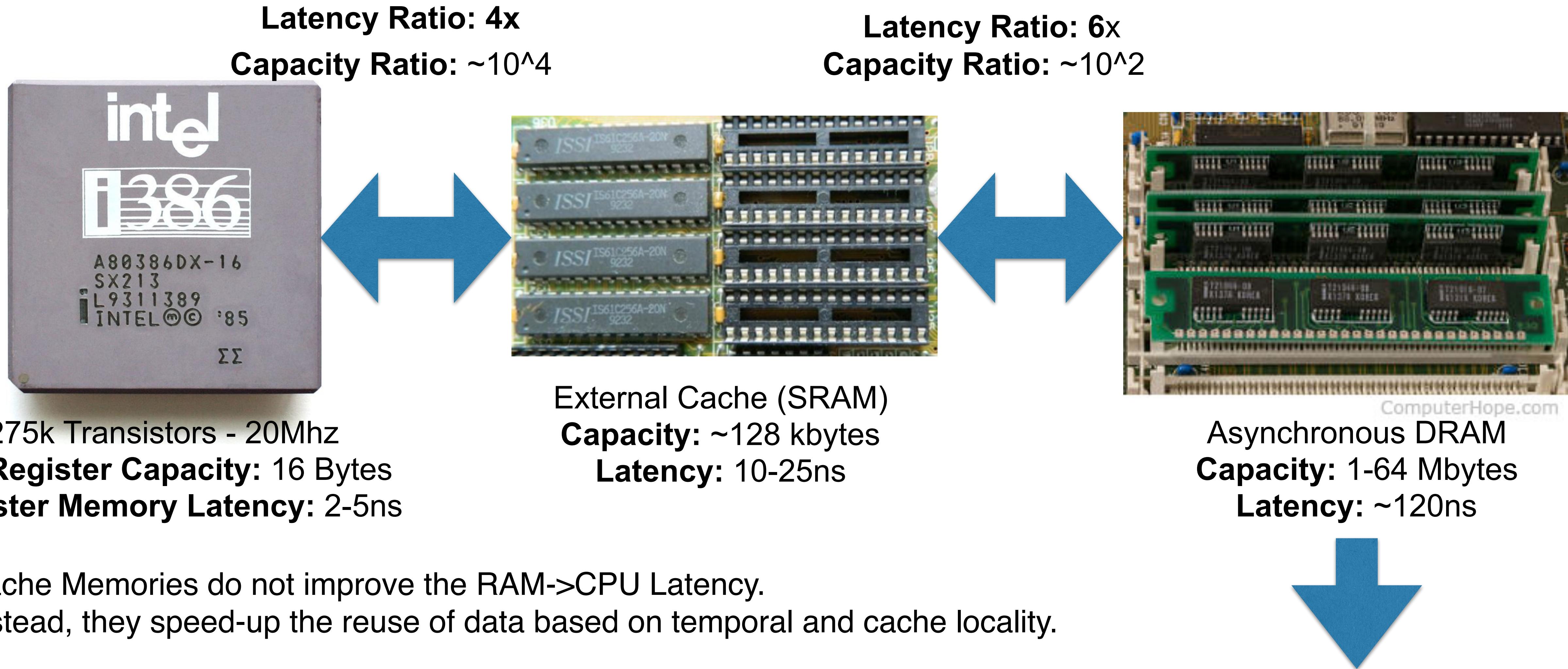
Problem - Growing disparity between register and RAM latencies.



© 2007 Elsevier, Inc. All rights reserved.

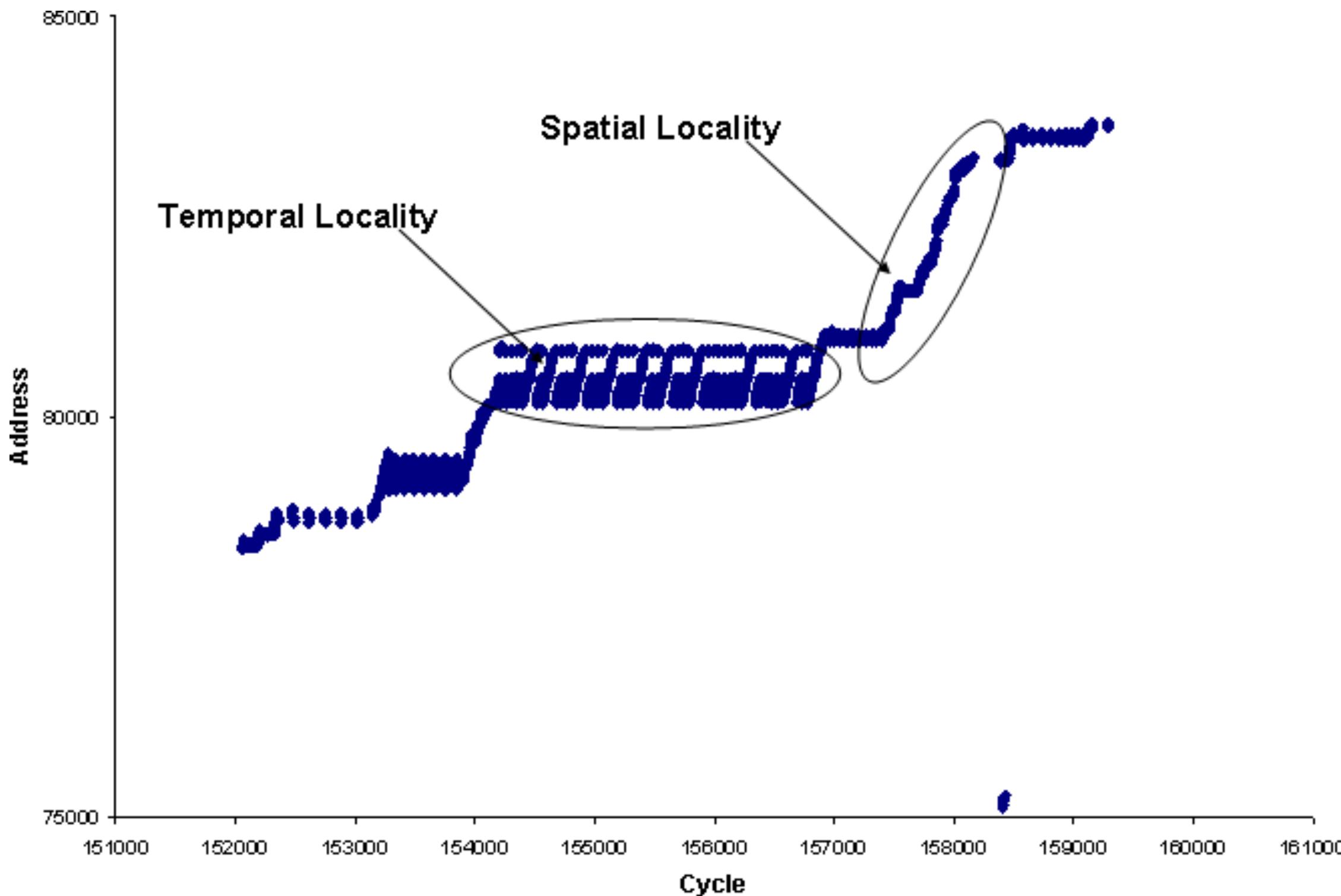
Cache Memory

1988 - Intel Releases 80386SX, the first commercial CPU with a Data-Cache Memory

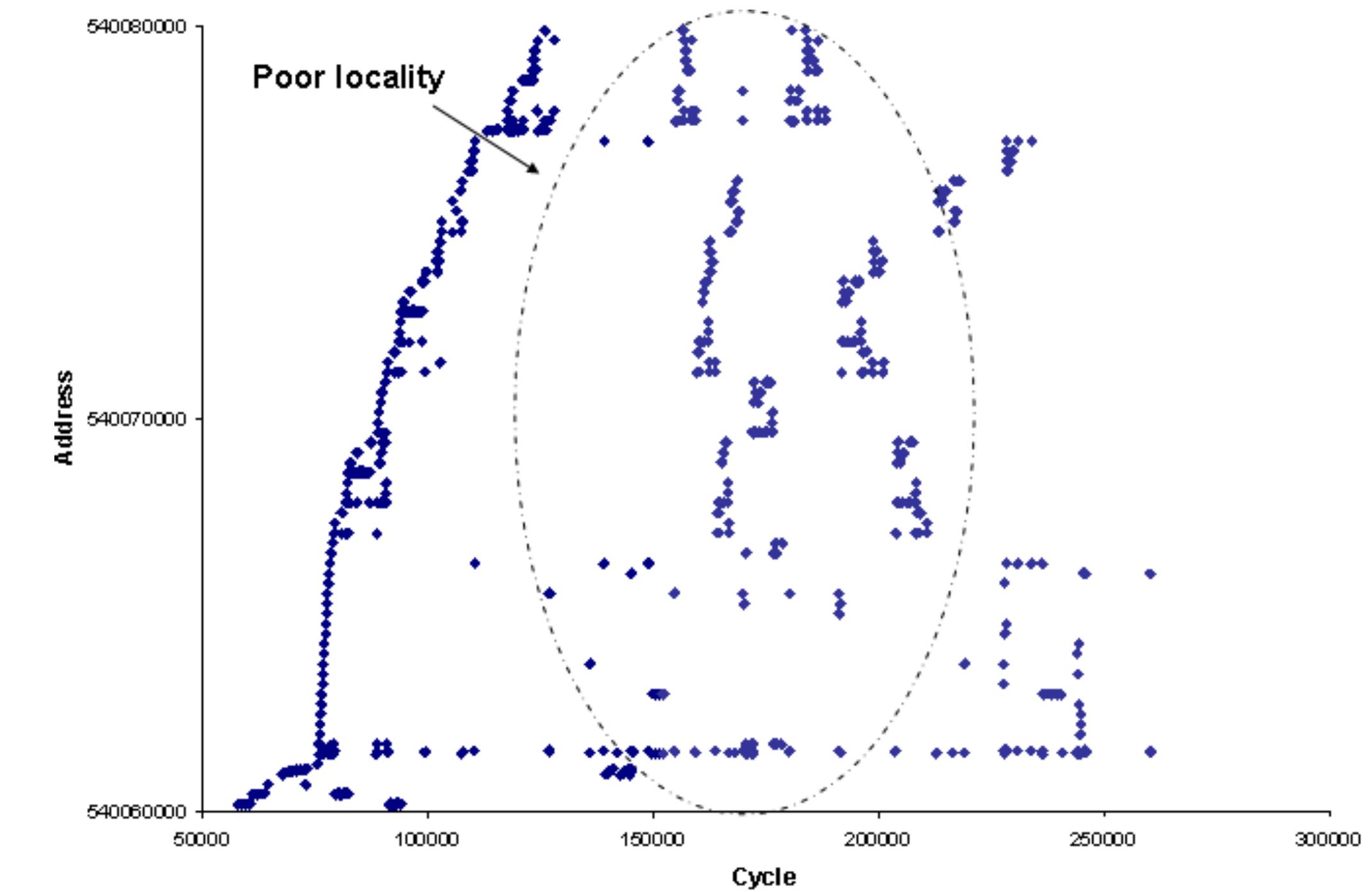


Cache Locality

Cache structures in modern processors benefit from both **Temporal** and **Spatial** locality.



High Cache Line Reuse



Frequent Cache Fails

Cache Hierarchy

201x - Modern Processors Employ Multiple Cache Levels

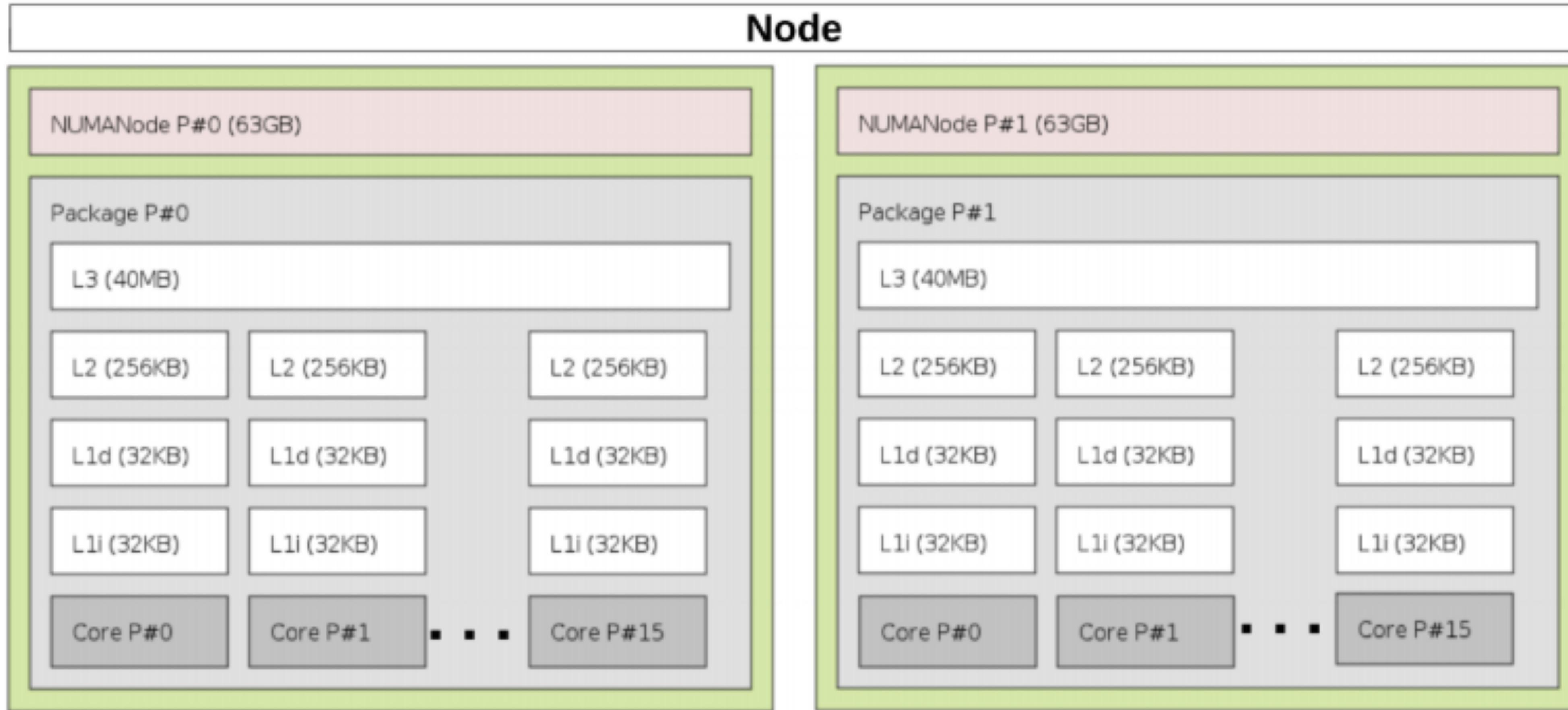
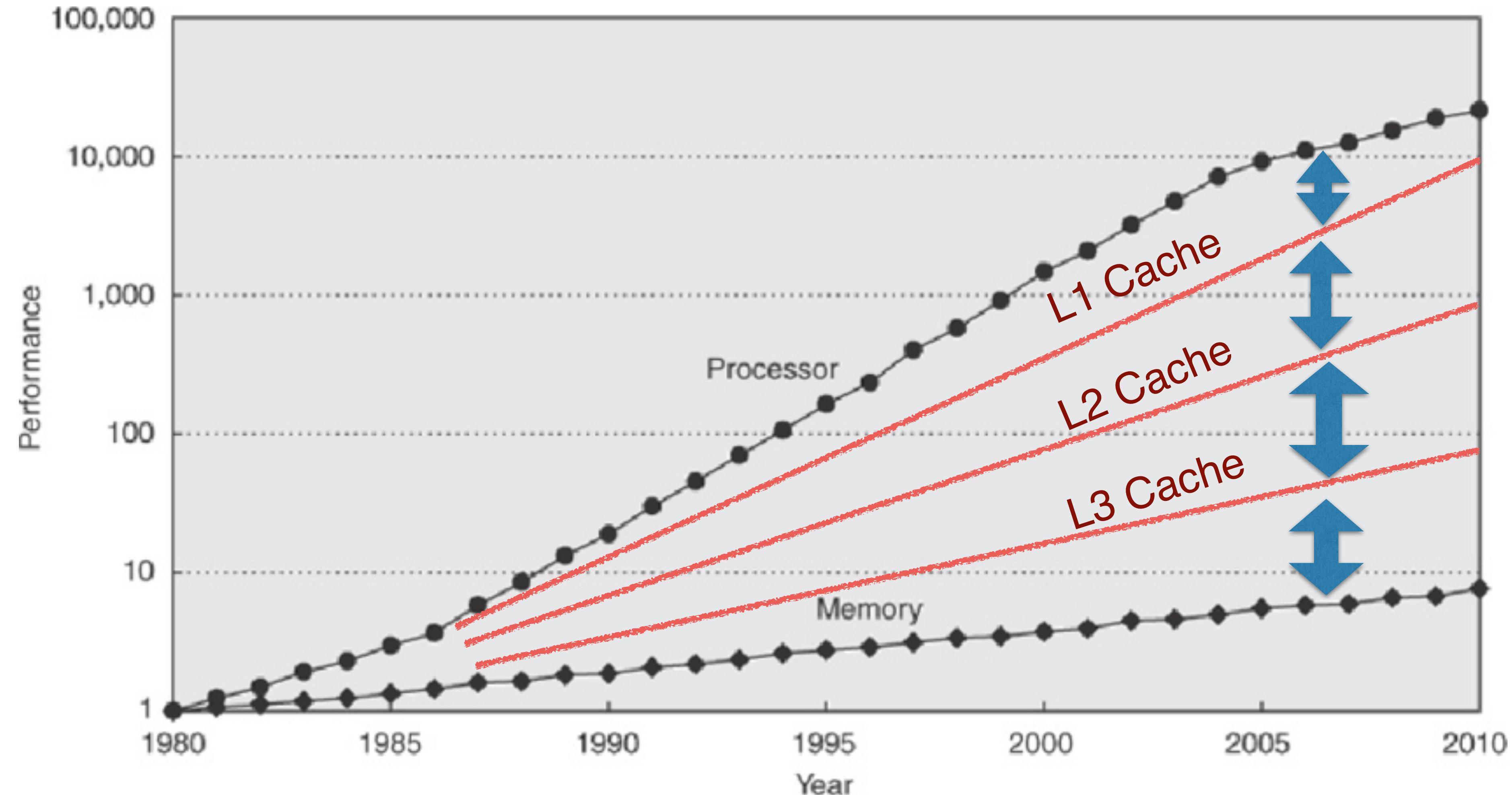


Figure 2.1: Configuration of a 32-core NERSC Cori Phase I (Haswell) Node.

RAM Latency

Solution - Cache memories bridges the gap between CPU and RAM performance

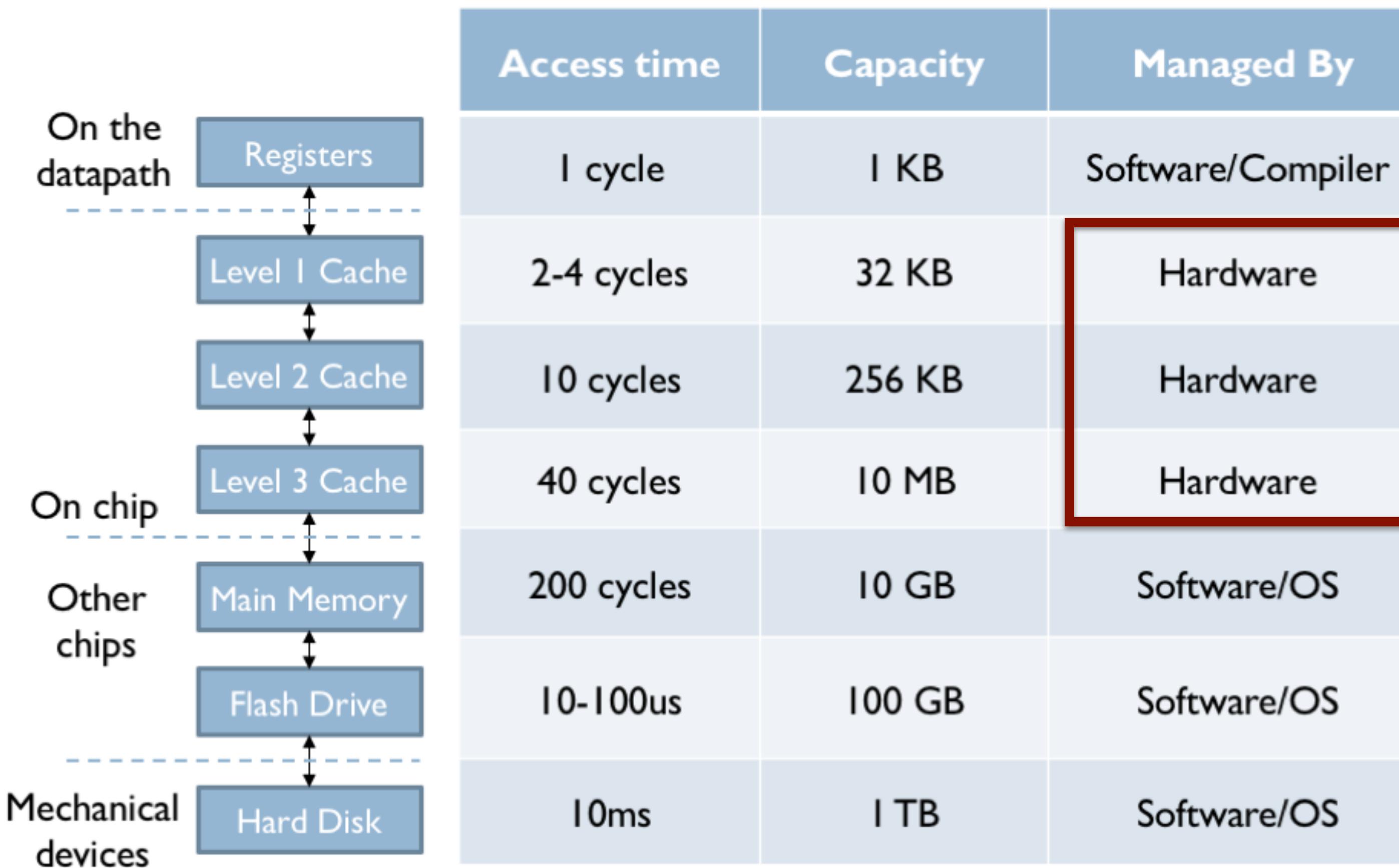


© 2007 Elsevier, Inc. All rights reserved.

Memory Hierarchy

A Typical Memory Hierarchy

- Everything is a cache for something else...

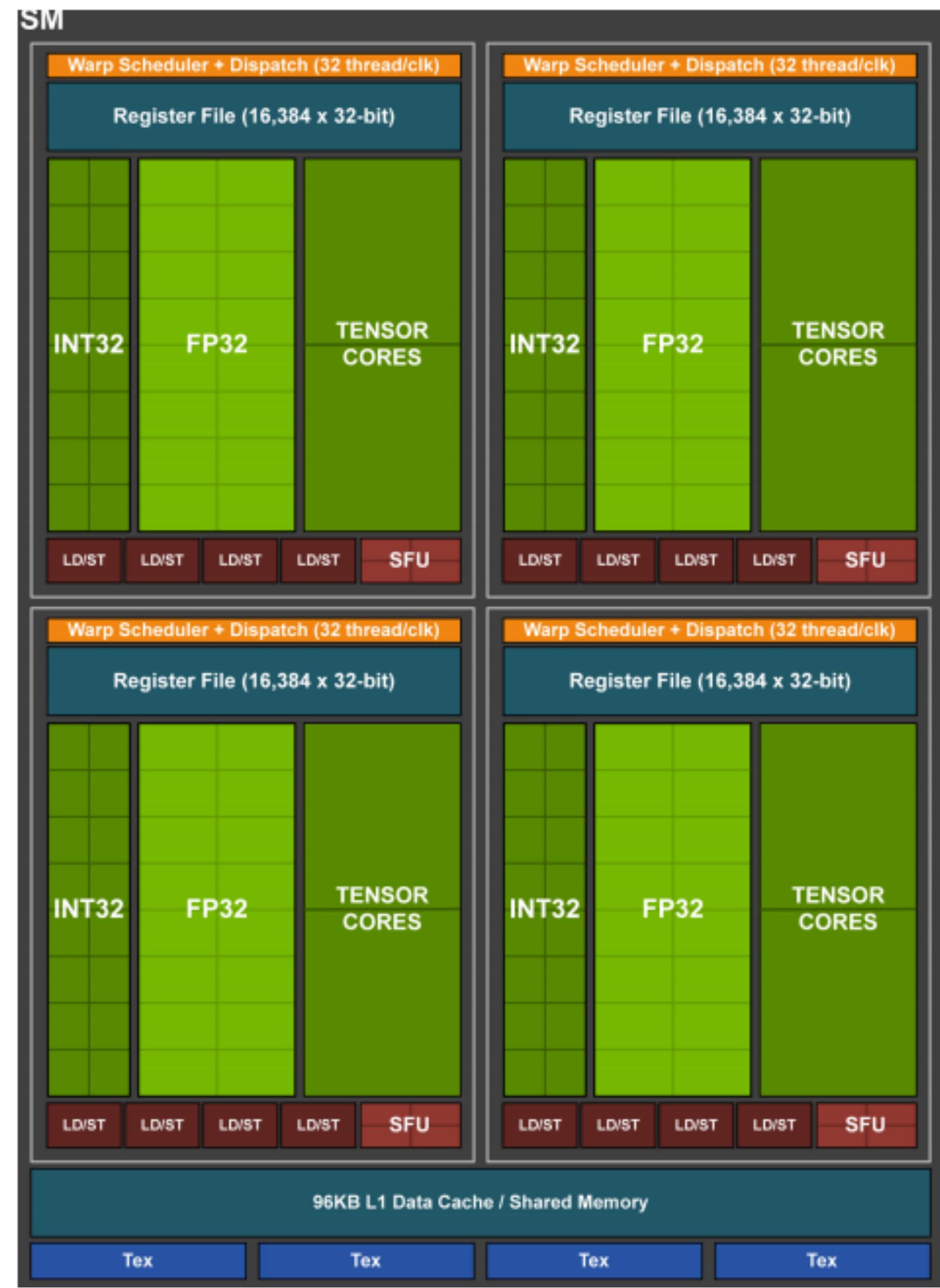
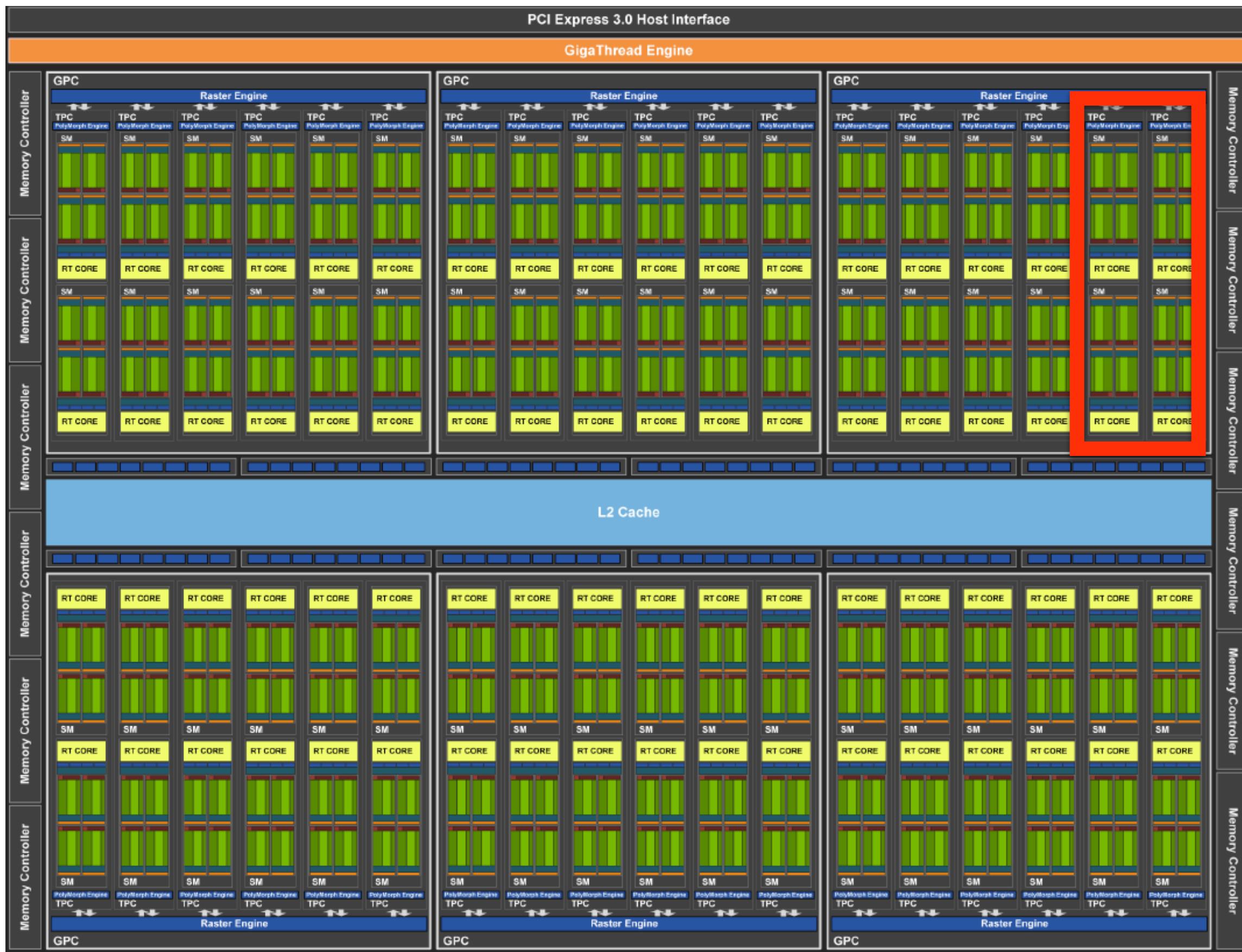


What does
this mean?

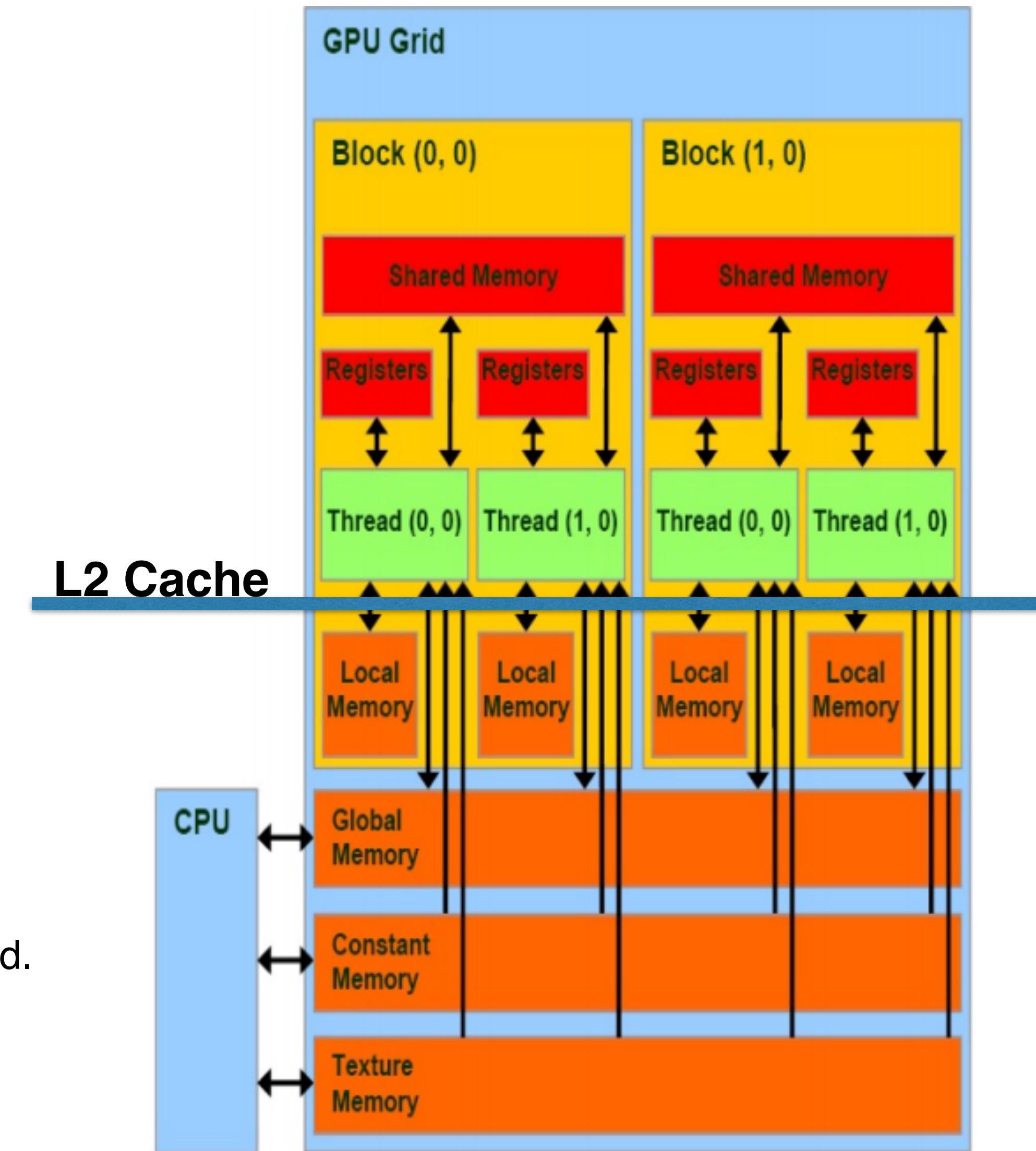
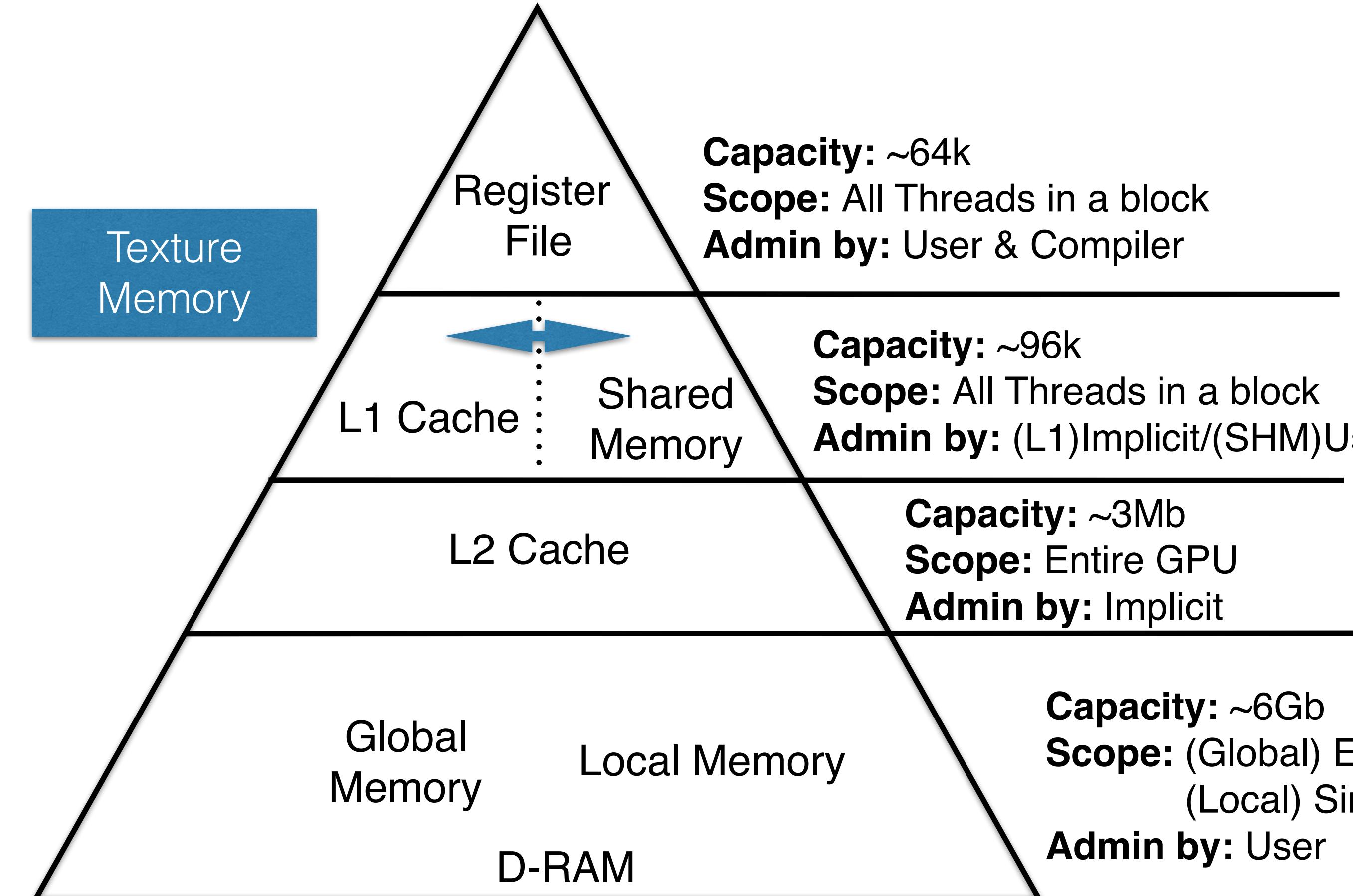
Memory Hierarchy on GPUs

Turing TU102 Streaming Multiprocessor(SM)

Up to 4608 threads executing concurrently.



GPU Memory Hierarchy



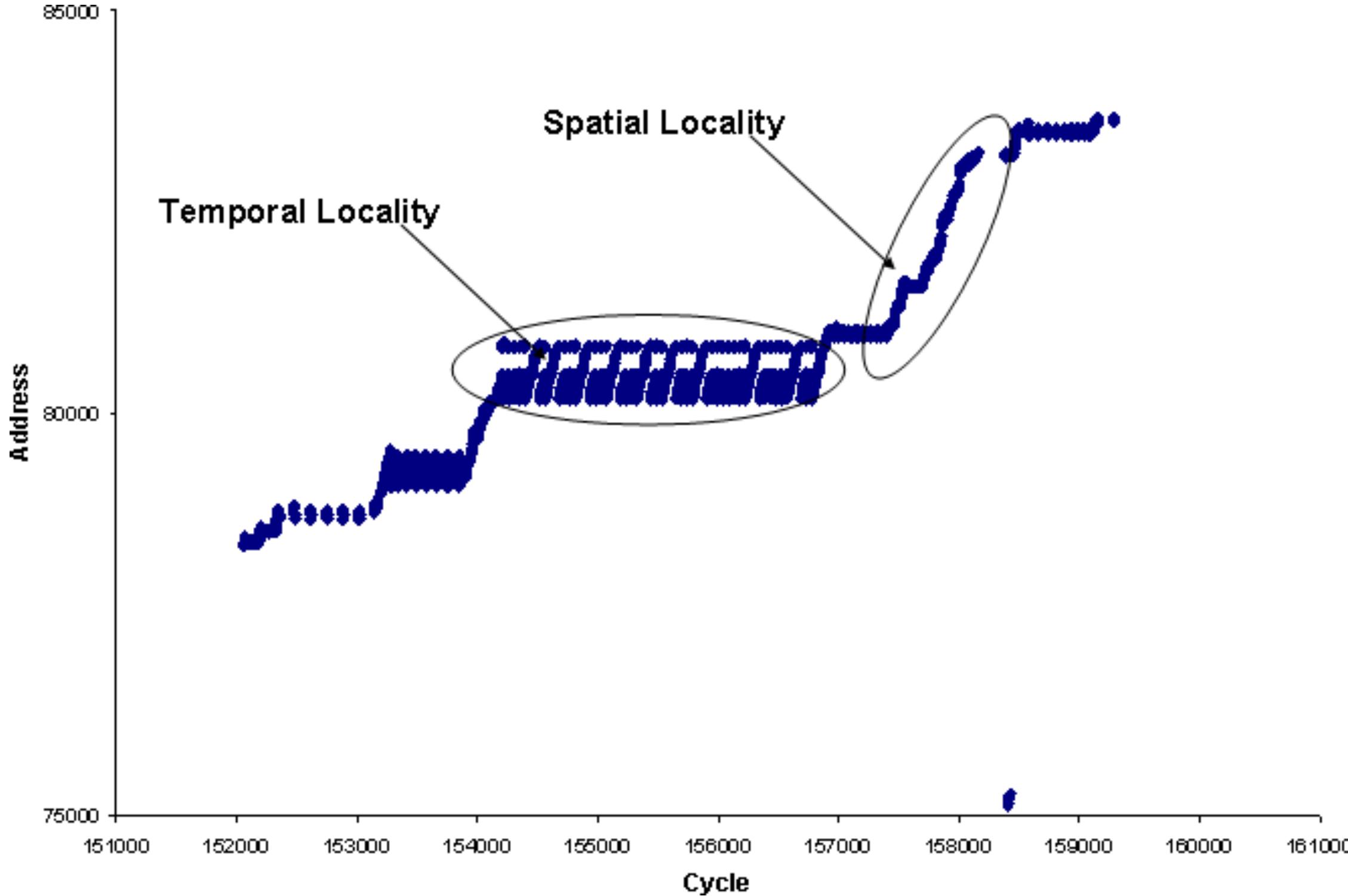
GPU Memory Hierarchy

Variable declaration	Memory	Scope	Lifetime
<code>__device__ __local__ int LocalVar;</code>	local	thread	thread
<code>__device__ __shared__ int SharedVar;</code>	shared	block	block
<code>__device__ int GlobalVar;</code>	global	grid	application
<code>__device__ __constant__ int ConstantVar;</code>	constant	grid	application

- `__device__` is optional when used with `__local__`, `__shared__`, or `__constant__`
- Automatic variables without any qualifier reside in a register
 - Except arrays that reside in local memory
 - Thread-local memory and spilled automatic variables is allocated in global memory

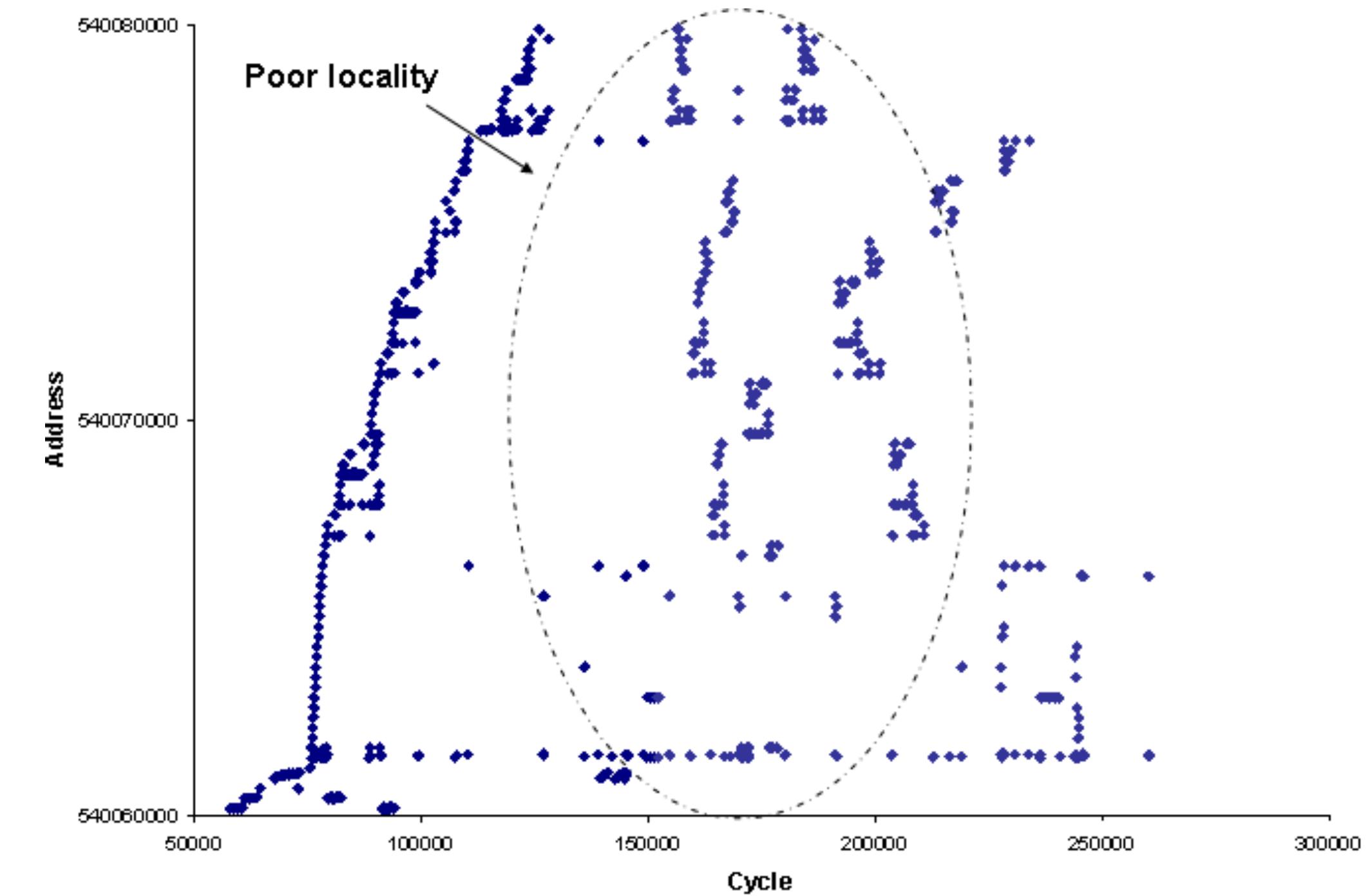
Cache Locality

Cache structures in modern processors benefit from both **Temporal** and **Spatial** locality.



High Cache Line Reuse

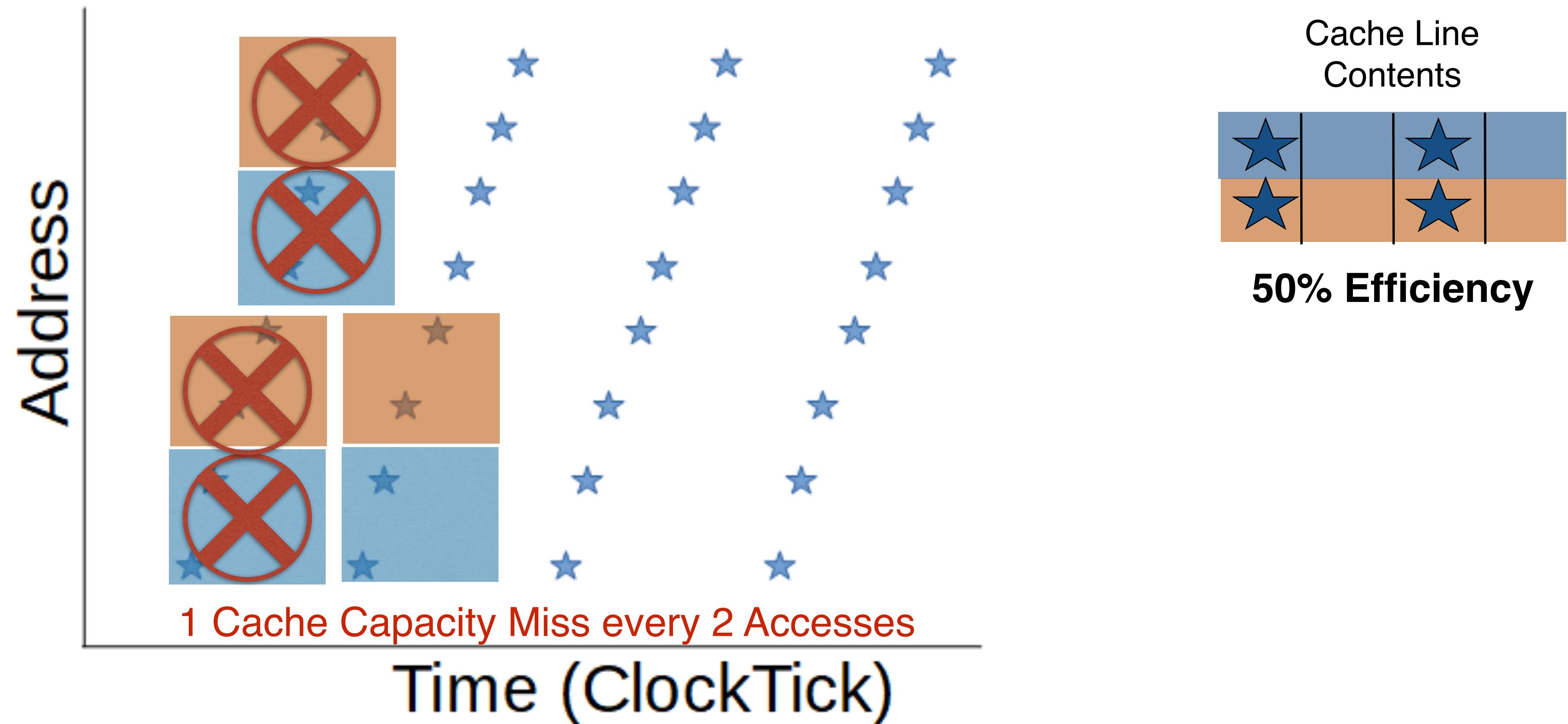
>> The caching/eviction decision is handled **implicitly** by the CPU guided by these heuristics <<
Dictates how a program has to be developed to harness its potential
This works well in the vast majority of cases, but may not always be the **optimal** strategy.



Frequent Cache Fails

Cache Locality

Consider the following case: CPU with 2 Cache Lines Available - 4 Elements per Line



Explicit Cache

The user decides what data is loaded into cache -> Manual Prefetching.

Pseudo-Example:

```
for (int i = 0; i < N; i++)  
cache[i] = prefetch(buf[i*2])
```



GPU Shared Memory

Key Points:

- Explicitly allocated/accessed by the programmer.
- Shared Among all threads in a block --> Requires proper synchronization.
- Remaining (unallocated) segments used for Local Memory and then L1 cache .
- Tailored to fit the application (not necessarily otherwise, like CPU cache)

Allocation:

- Host code: On Kernel Launch

```
size_t shmSize = 1000*sizeof(double); // The rest goes to L1 Cache  
myKernel<<<grid, block, shmSize>>>(pars);
```

Shared Memory Size should not exceed the L1-cache size of your GPU

Usage:

- Device code: Inside the kernel

```
extern __shared__ double mySharedArray[];  
s[i] = globalArray[i*n]; // Prefetching global memory into shared memory
```

Register / Local Memory

Register Optimization

- Registers are explicitly allocated/accessed by the programmer.

CPU:

Predefined by Hardware
through register names
Employed by the compiler
RAX ... RBX ...RCX

GPU:

General Space (64kb) for entire Block
Programmer defines and names registers:
double mySum = 0.0;
double myAverage = 0.0;

- GPU Compiler is less aggressive with optimizations due assumed parallelism
- Limited number of registers:
NRegs * NThreads should not exceed Register File Size (e.g., 64kb)

Warning: Excessive Registers will spill to local memory (DRAM)

- Exceptions:
 - Pre-defined arrays (e.g., myArray[512]) are stored in local memory.
 - Constant variables (e.g., const myPI = 3.14) may not count (are optimized out)

GPU Shared Memory (Static)

Example: 2D Problem - 4k Elements per side

```
#define BLOCKSIZE 32
N = 4096;
dim3 blockDim(BLOCKSIZE,BLOCKSIZE) // 512 Threads
dim3 gridDim(n/BLOCKSIZE, n/BLOCKSIZE) // 128x128 Blocks = 16k Blocks
my2dProblem<<<gridDim, blockDim>>>(grid1, grid2, N);
```

Device Code:

```
__global__ void my2DProblem(double* grid1, double* grid2, size_t N)
{
    const size_t m      = blockIdx.x * blockDim.x + threadIdx.x;
    const size_t t      = threadIdx.x;

    __shared__ double s0[BLOCKSIZE][BLOCKSIZE]; // 32*32*8 bytes = 8k bytes
    __shared__ double s1[BLOCKSIZE][BLOCKSIZE]; // 32*32*8 bytes = 8k bytes

    s[2*t + 0] = grid1[m];
    s[2*t + 1] = grid2[m];

    __syncthreads(); // Important: synchronize to make sure every thread in the block saved their values before accessing.
    double myCalc = 0.0;
    for (size_t i = 0; i < blockDim.x; i++) myCalc += sqrt(s0[i] + s1[i]);
}
```

GPU Shared Memory (Dynamic)

Example: 2D Problem - 4k Elements per side

```
N = 4096;  
dim3 blockDim(32,32) // 512 Threads  
dim3 gridDim(n/32, n/32) // 128x128 Blocks = 16k Blocks  
size_t shmSize = 2*32*32*sizeof(double);  
my2dProblem<<<gridDim, blockDim, shmSize>>>(grid1, grid2, N);
```

Shared Memory:

2 Doubles per Thread = $512 \times 2 \times 8\text{bytes} = 8\text{kb}$

Device Code:

```
__global__ void my2DProblem(double* grid1, double* grid2, size_t N)  
{  
    const size_t m      = blockIdx.x * blockDim.x + threadIdx.x;  
    const size_t t      = threadIdx.x;  
  
extern __shared__ double s[];  
s[2*t + 0] = grid1[m];  
s[2*t + 1] = grid2[m];  
  
__syncthreads(); // Important: synchronize to make sure every thread in the block saved their values before accessing.  
double myCalc = 0.0;  
for (size_t i = 0; i < blockDim.x; i++) myCalc += sqrt(s[2*i] + s[2*i+1]);  
}
```

Register Memory:

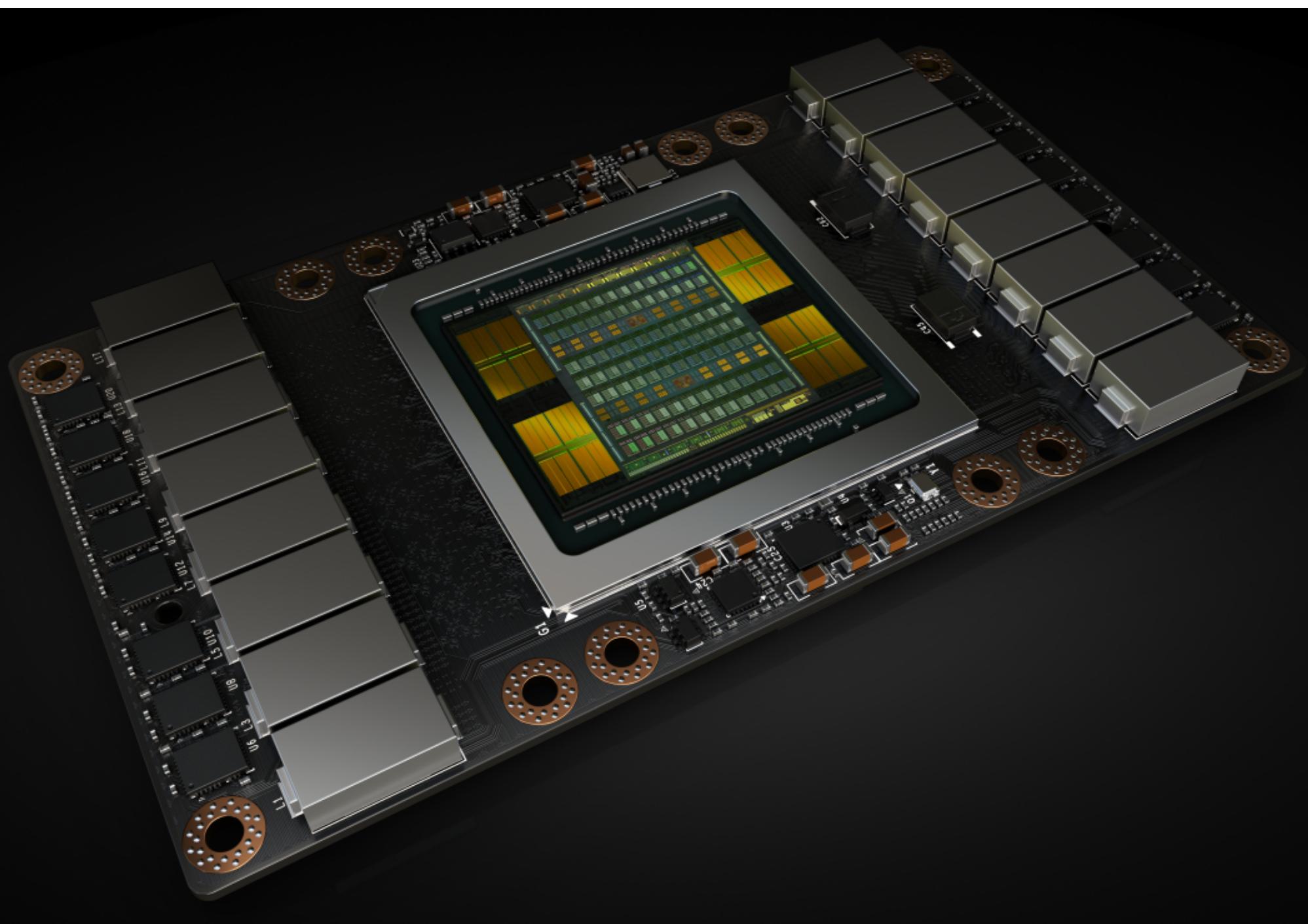
7 Registers * 512 Threads * 8bytes = 28kb

No spillage to local memory.

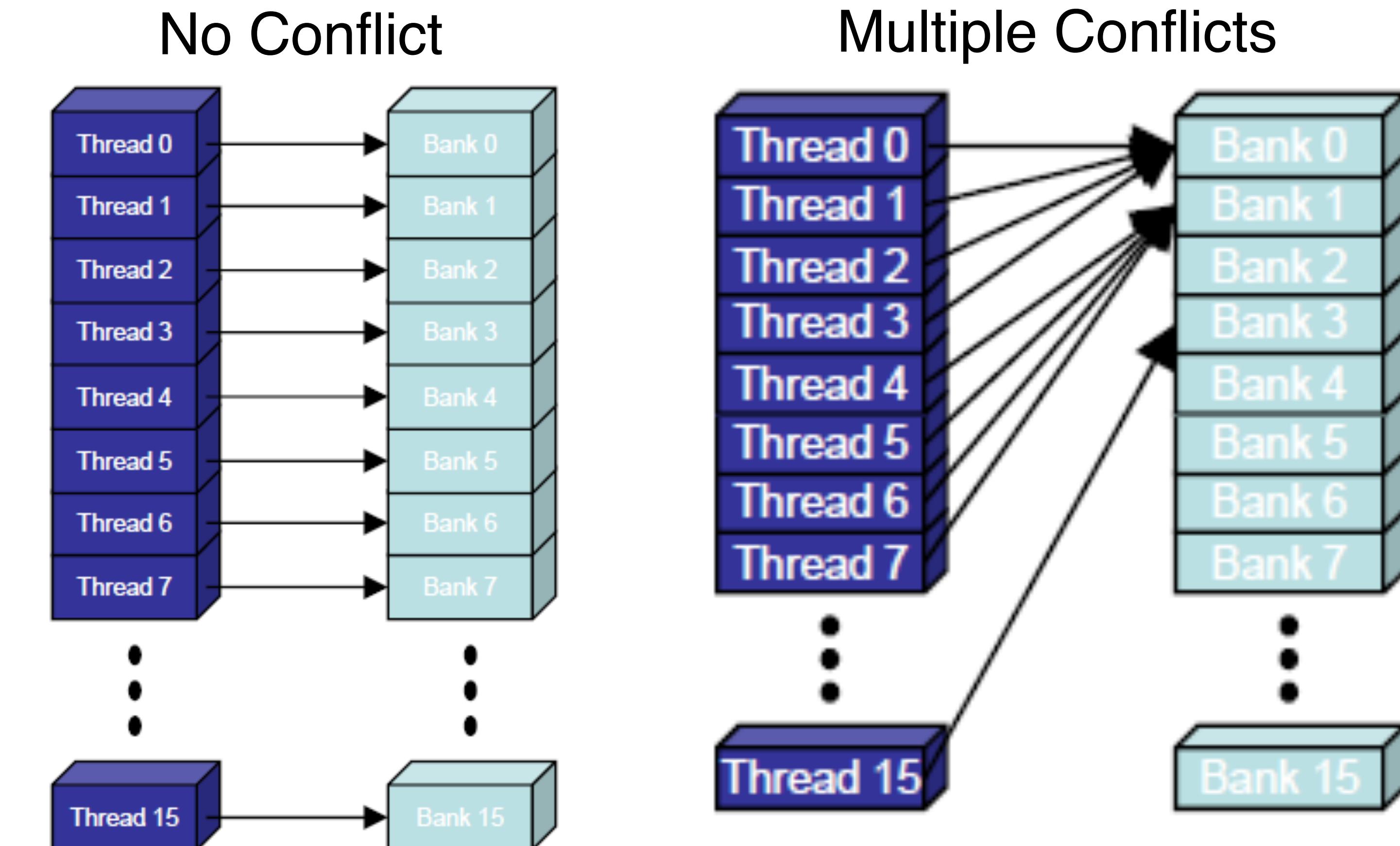
Banked Access to Global & Shared Memory

GPU RAM

Optimized for High-bandwidth
Multi-banked access.



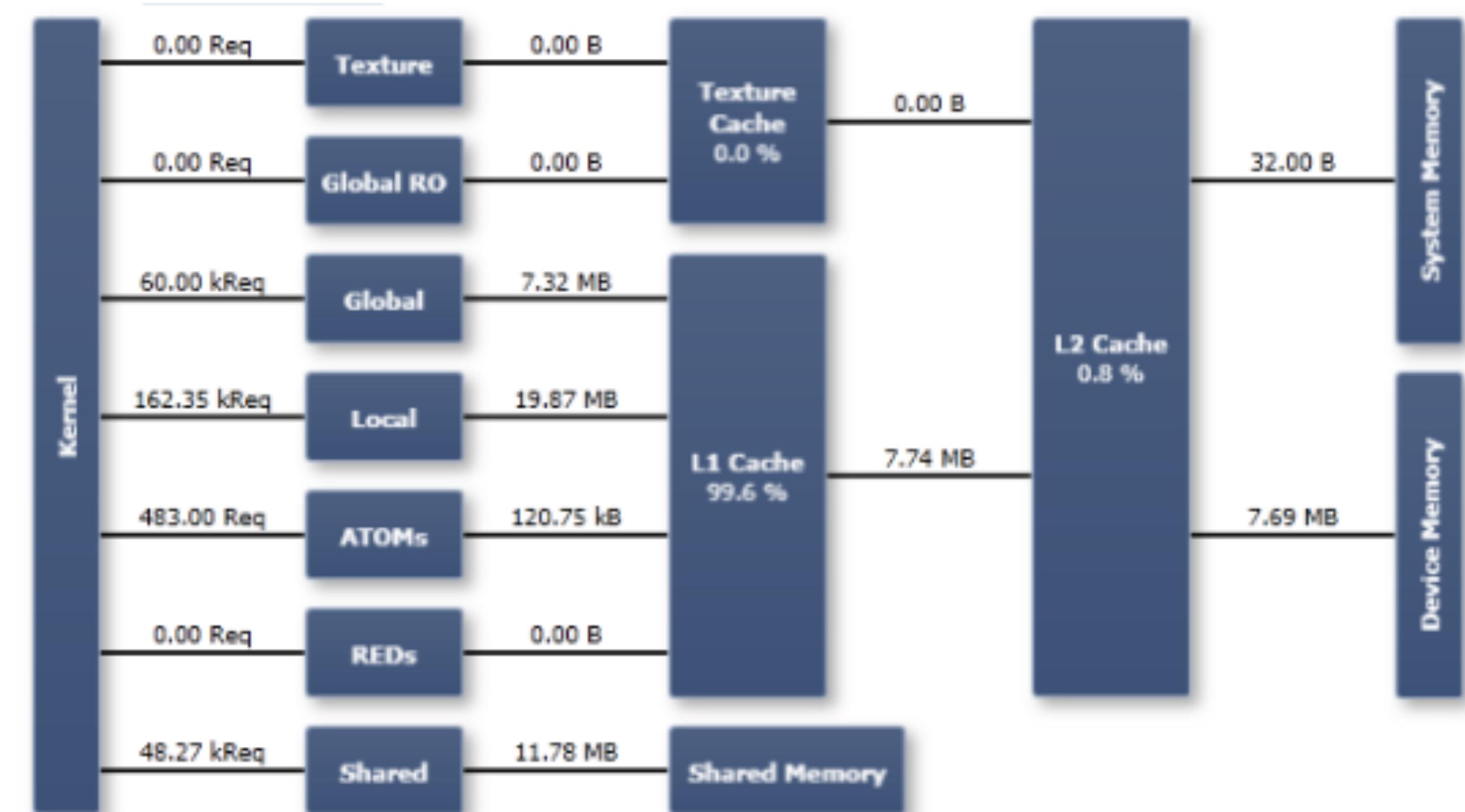
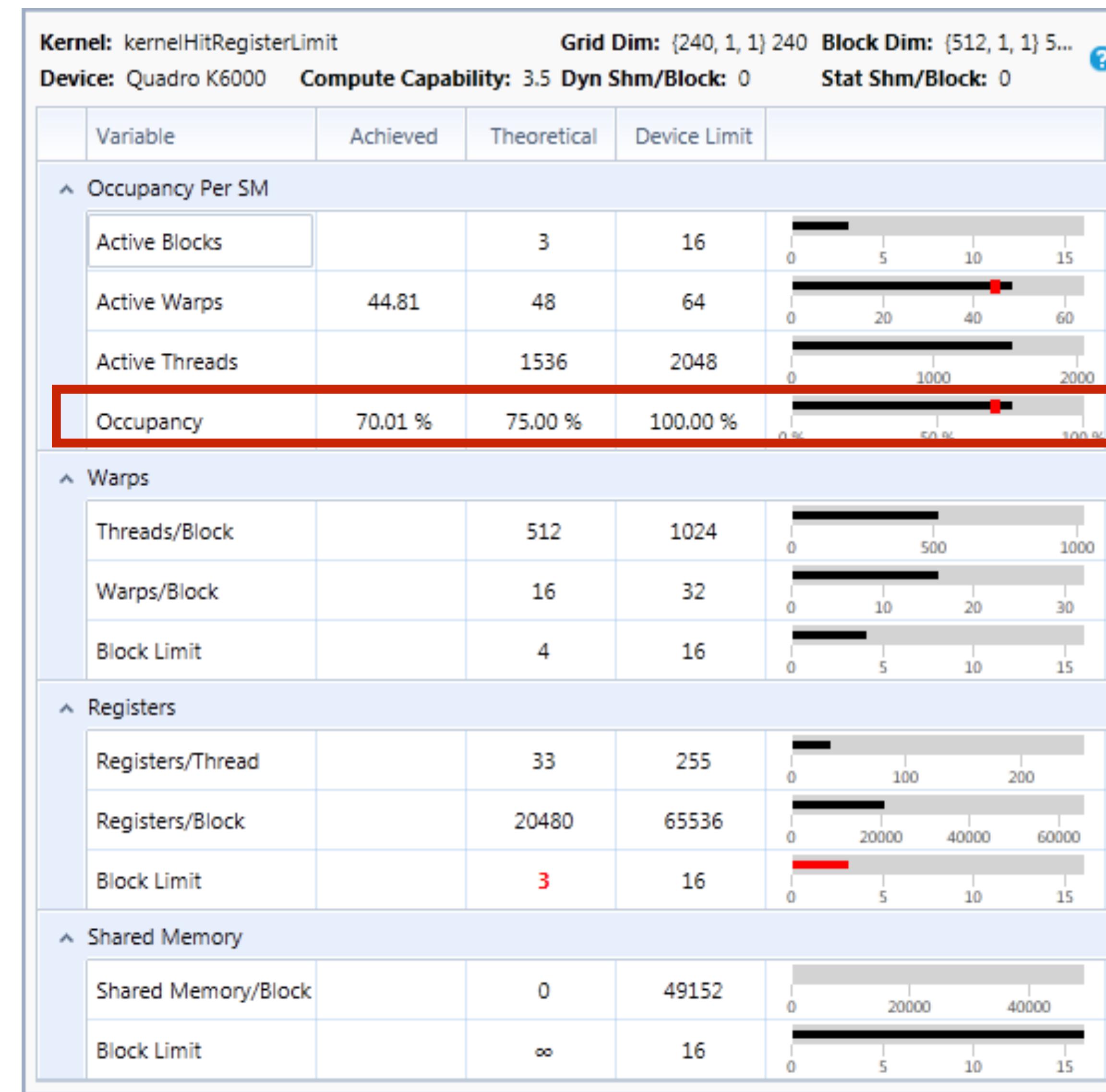
Source: [NVIDIA](#)



Source: [CUDA Programming Blog](#)

Profiling

NVIDIA Nsight - Displays the achieved (compute) occupancy and memory behavior of each kernel.



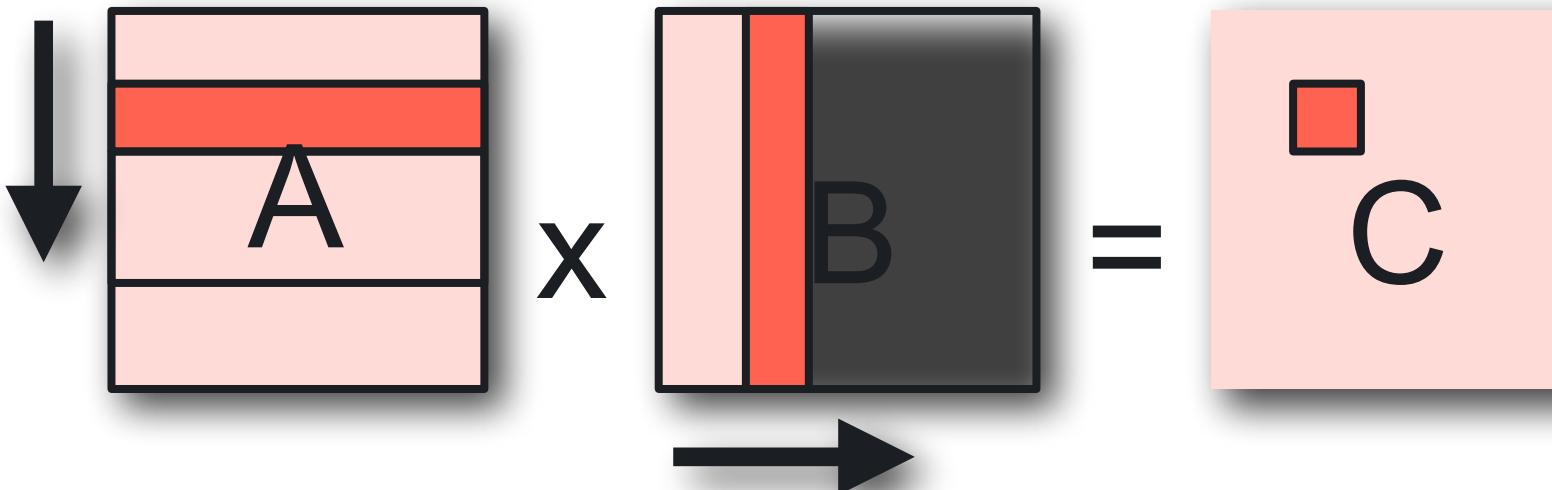
Uses:

- Squeeze computational capacity of the GPU
- Measure efficiency of Cache/SHM optimizations

Example

Matrix/Matrix Multiplication

Square Matrix Multiplication

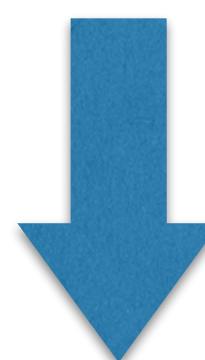


```
for (int i=0; i<N; i++)  
    for (int j=0; j<N; j++)  
        for (int k=0; k<N; k++)  
            c[i*N+j] += a[i*N+k] * b[k*N+j];
```

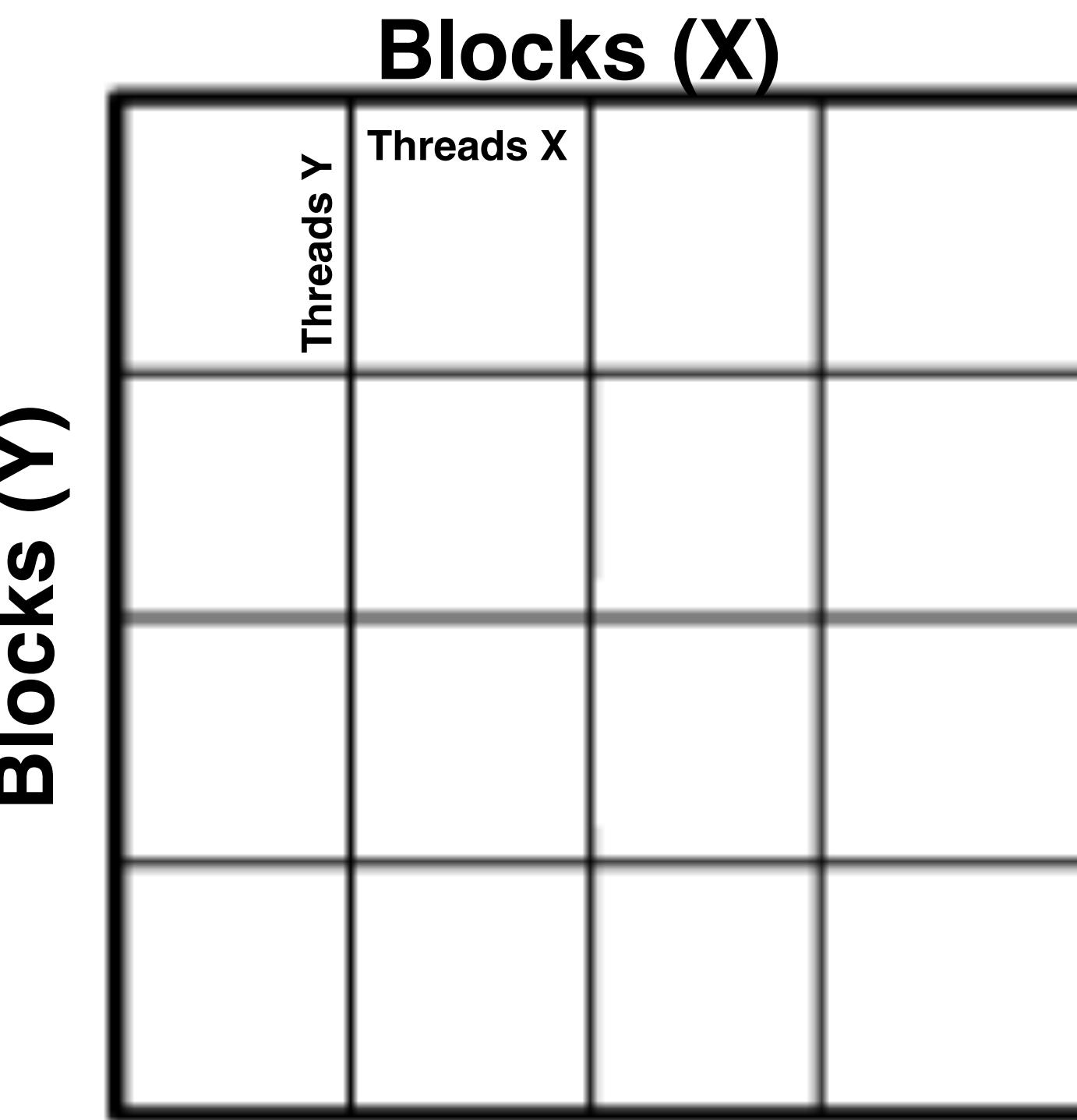
2D Matrices (NxN)
1D (row by column) operation
Complexity: $O(N^3)$

Parallelize
2D Geometry

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



Blocks (X) Threads (X)
Blocks (Y) Threads (Y)



Each thread
Executes a 1D Kernel

```
for (int k=0; k<N; k++)  
    c[i*N+j] += a[i*N+k] * b[k*N+j];
```

Naive matrix multiplication kernel

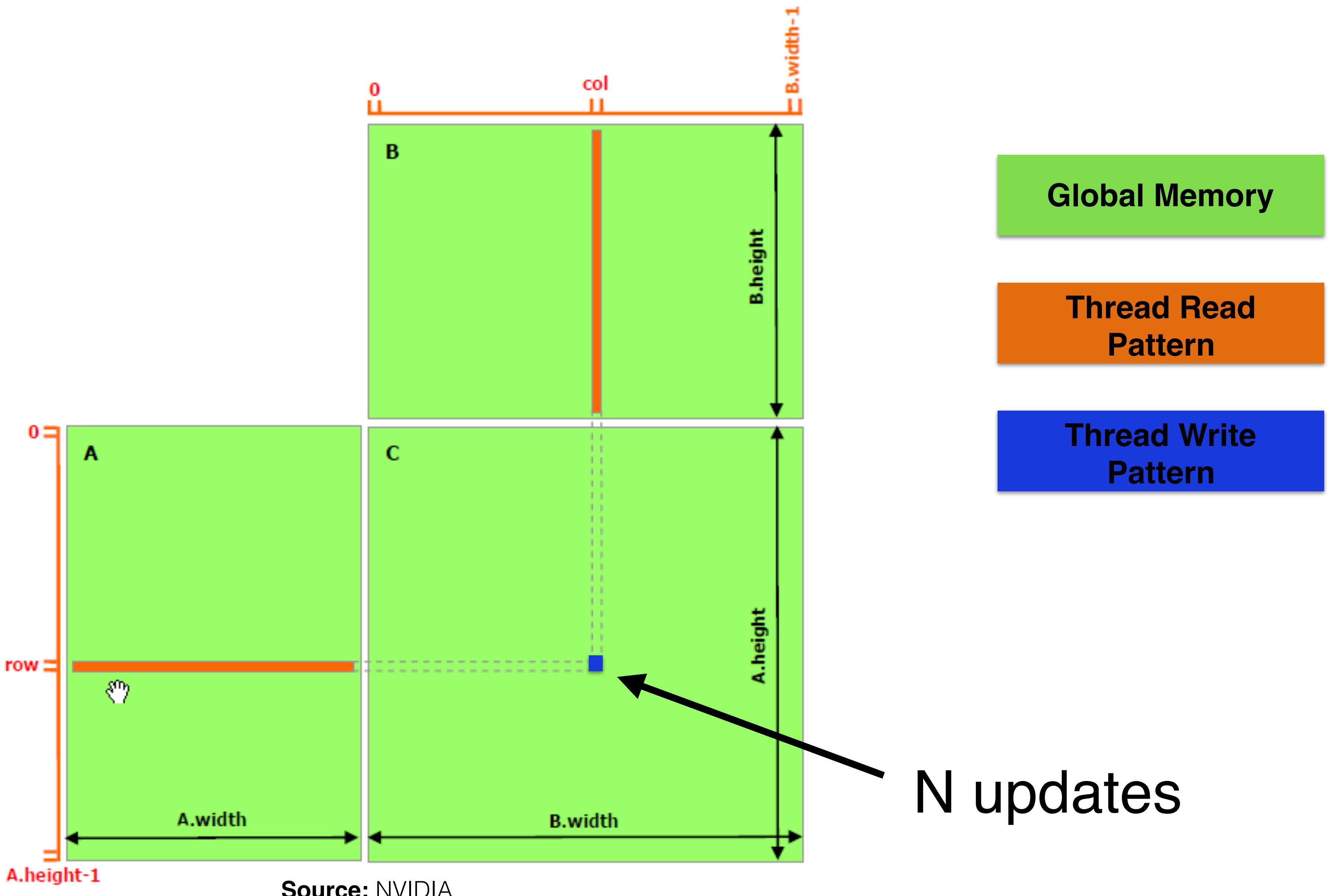
```
// Host Code
#define BLOCK_SIZE 32
dim3 blocks(N/BLOCK_SIZE, N/BLOCK_SIZE);
dim3 threads(BLOCK_SIZE, BLOCK_SIZE); // 512 Threads per Block
matrix<<< blocks, threads >>>(dev_a, dev_b, dev_c, N);

// Device Code
__global__ void matrix(double* A, double* B, double* C, int N)
{
    int x = threadIdx.x + blockIdx.x*blockDim.x;
    int y = threadIdx.y + blockIdx.y*blockDim.y;

    if (x<N && y<N)
    {
        c[x*N + y] = 0;

        for (int k=0; k<N; k++)
            C[x*N + y] += A[x*N + k] * B[k*N + y];
    }
}
```

Memory Access Pattern



Shared Memory Optimized Kernel

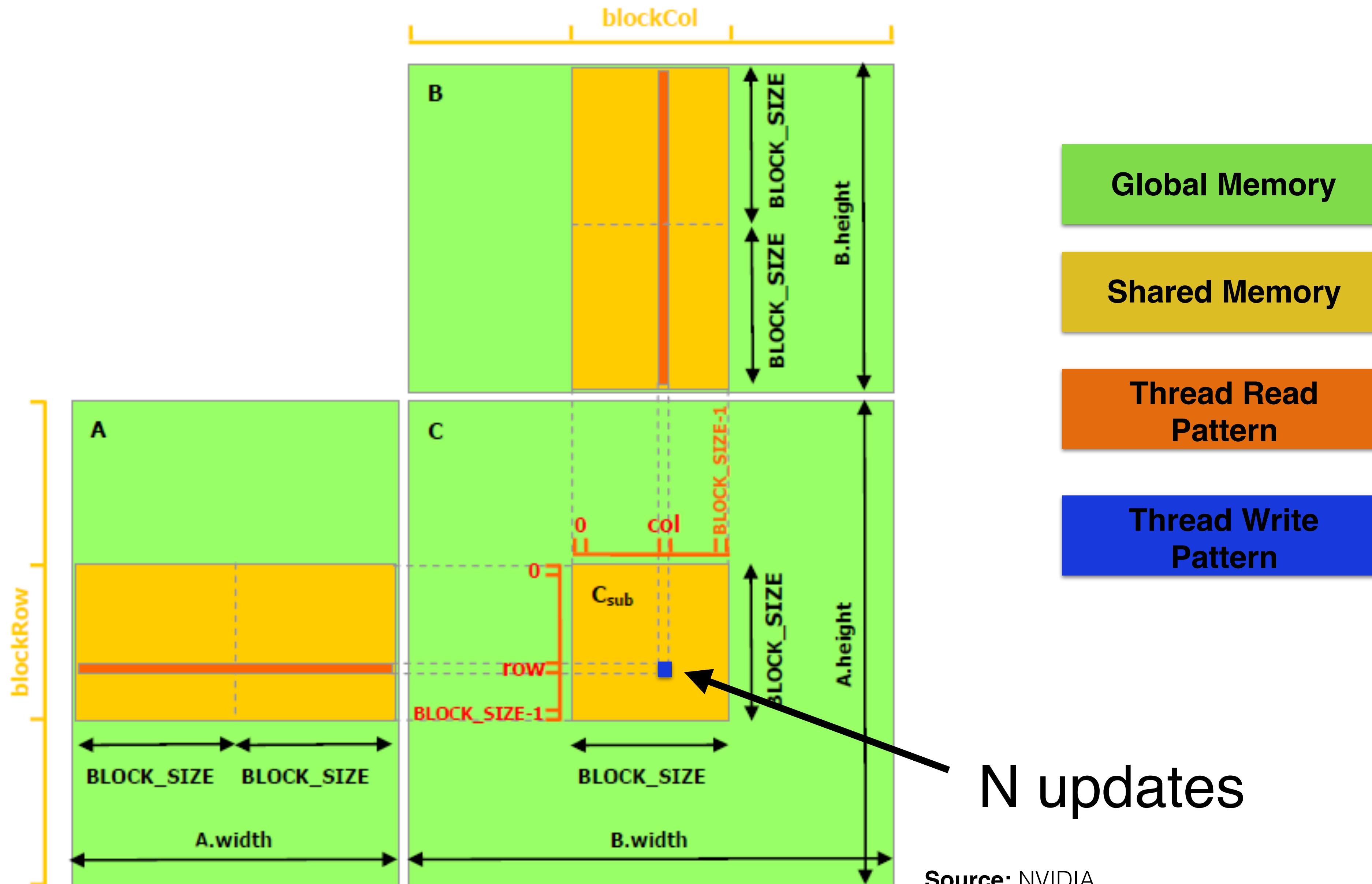
```
__global__ void matrix(double* A, double* B, double* C, int N)
{
    int x = threadIdx.x + blockIdx.x*blockDim.x;
    int y = threadIdx.y + blockIdx.y*blockDim.y;
    int myRow = threadIdx.y;
    int myCol = threadIdx.x;

    for (int m = 0; m < (N / BLOCK_SIZE); m++)
    {
        __shared__ double As[BLOCK_SIZE][BLOCK_SIZE];
        __shared__ double Bs[BLOCK_SIZE][BLOCK_SIZE];

        As[myRow][myCol] = A[x*N + m*BLOCK_SIZE+myCol];
        Bs[myRow][myCol] = B[(m*BLOCK_SIZE + myRow)*N + y];

        __syncthreads();
        for (int i = 0; i < BLOCK_SIZE; i++) C[x*N + y] += As[myRow][i] * Bs[i][myCol];
        __syncthreads();
    }
}
```

Memory Access Pattern



One more optimization:

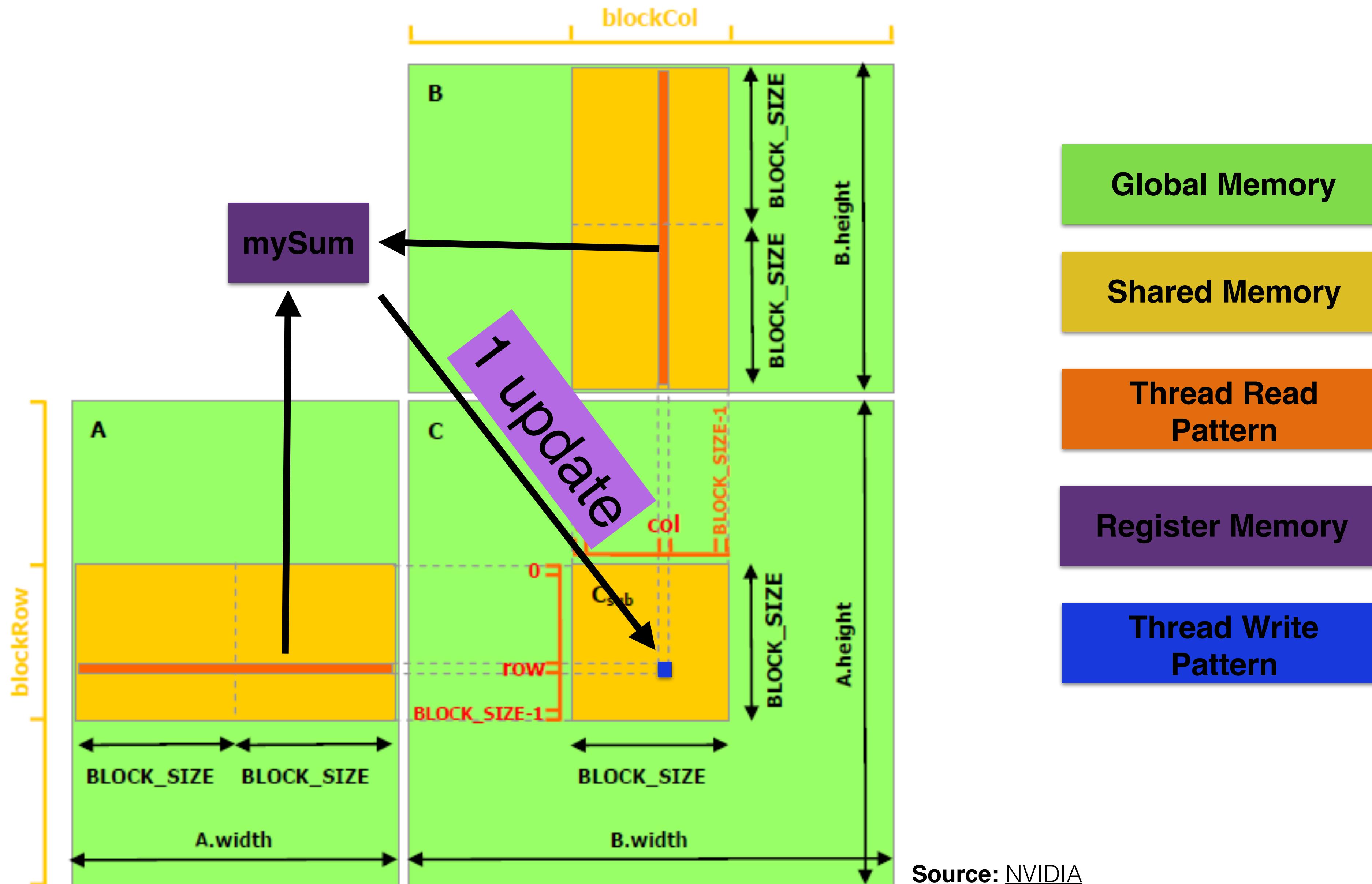
```
__global__ void matrix(double* A, double* B, double* C, int N)
{
    int x = threadIdx.x + blockIdx.x*blockDim.x;
    int y = threadIdx.y + blockIdx.y*blockDim.y;
    int myRow = threadIdx.y
    int myCol = threadIdx.x

    double mySum = 0;
    for (int m = 0; m < (N / BLOCK_SIZE) ; m++)
    {
        __shared__ double As[BLOCK_SIZE][BLOCK_SIZE];
        __shared__ double Bs[BLOCK_SIZE][BLOCK_SIZE];

        As[myRow][myCol] = A[x*N + m*BLOCK_SIZE+myCol];
        Bs[myRow][myCol] = B[(m*BLOCK_SIZE + myRow)*N + y];

        __syncthreads();
        for (int i = 0; i < BLOCK_SIZE; i++) mySum += As[myRow][i] * Bs[i][myCol];
        __syncthreads();
    }
    C[x*N + y] += mySum;
}
```

Memory Access Pattern

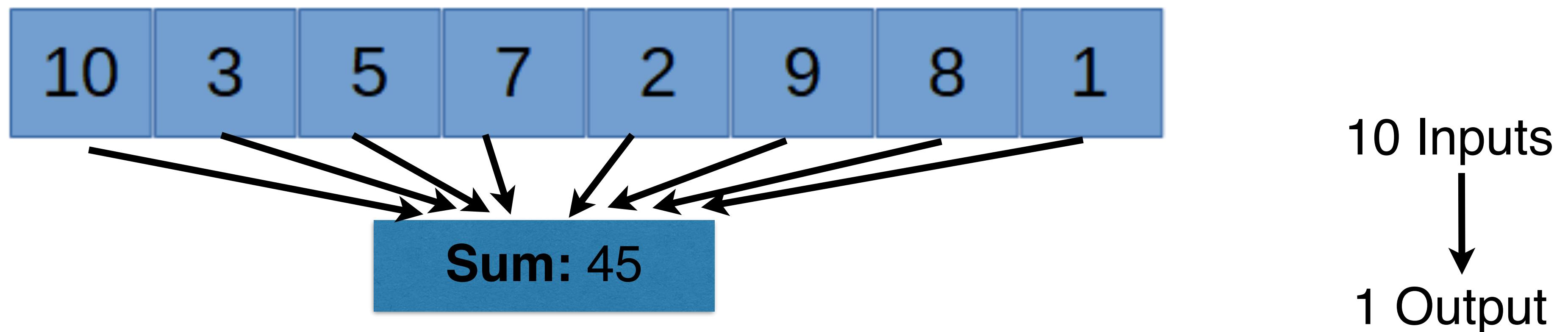


Kernel Decomposition

Vector Reduction

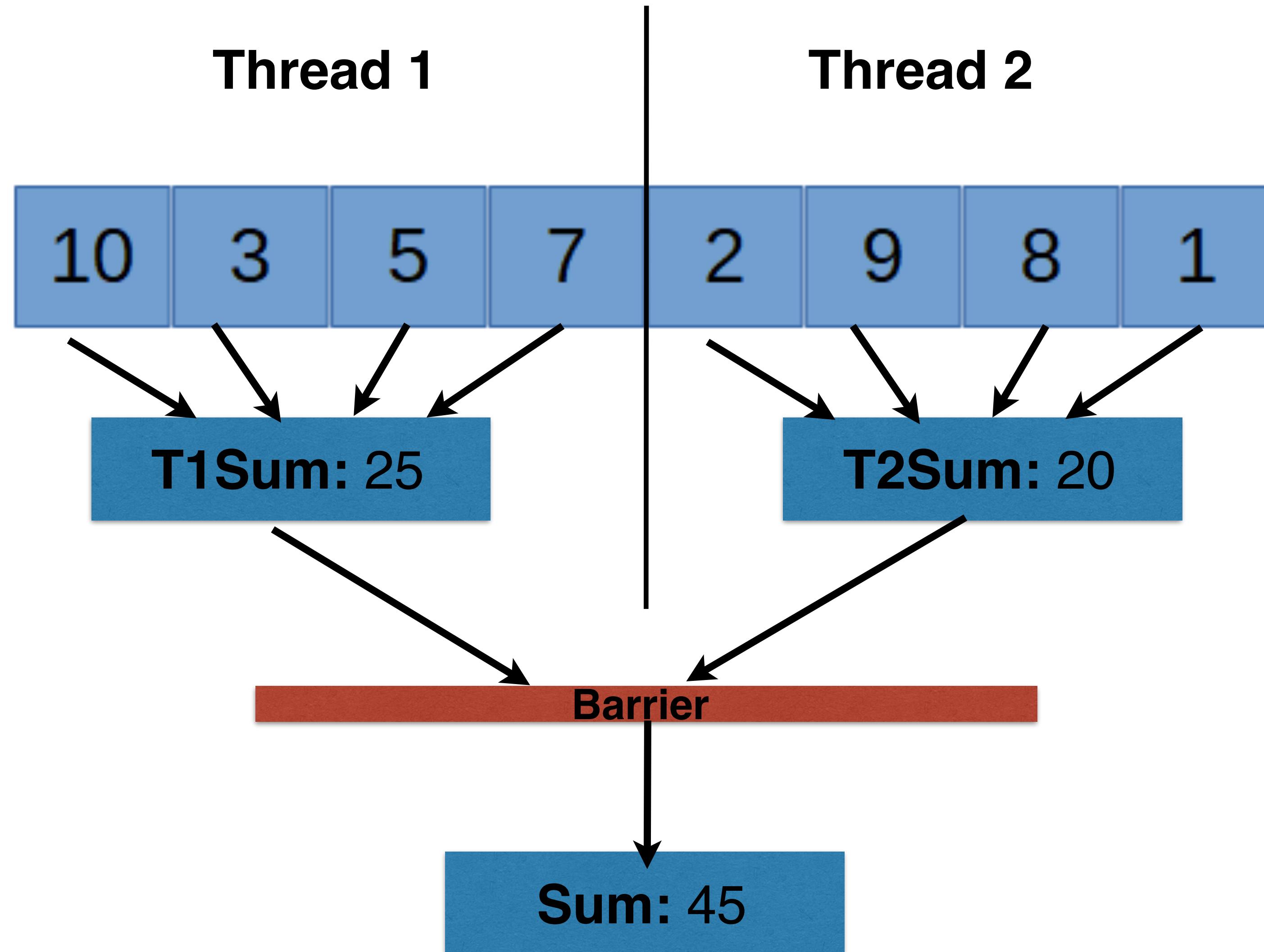
Fact: Some problems require reading from an entire array but write to a smaller set.

Example: Vector reduction



```
double mySum = 0.0;  
for (int i = 0; i < N; i++) mySum += vector[i];
```

Parallel Vector Reduction



Distributed Parallel Vector Reduction

