

Optimizing Algorithms and Code for Data Locality and Parallelism

Marc Moreno Maza

University of Western Ontario, Canada

SHARCNET Summer Seminar, London
June 17, 2013

What is this tutorial about?

Optimizing algorithms and code

- Improving code performance is hard and complex.
- Requires a good understanding of the underlying algorithm and implementation environment (hardware, OS, compiler, etc.).

What is this tutorial about?

Optimizing algorithms and code

- Improving code performance is hard and complex.
- Requires a good understanding of the underlying algorithm and implementation environment (hardware, OS, compiler, etc.).

Optimizing for data locality

- Computer cache memories have led to introduce a new complexity measure for algorithms and new performance counters for code.
- Optimizing for data locality brings large speedup factors.

What is this tutorial about?

Optimizing algorithms and code

- Improving code performance is hard and complex.
- Requires a good understanding of the underlying algorithm and implementation environment (hardware, OS, compiler, etc.).

Optimizing for data locality

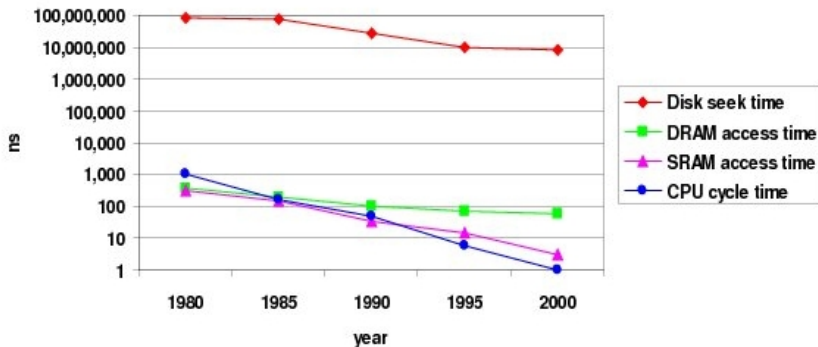
- Computer cache memories have led to introduce a new complexity measure for algorithms and new performance counters for code.
- Optimizing for data locality brings large speedup factors.

Optimizing for parallelism

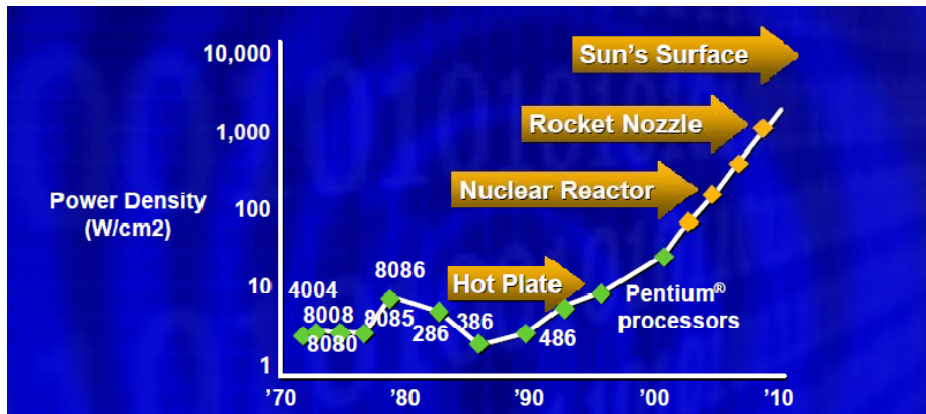
- All recent home and office desktops/laptops are parallel machines; moreover “GPU cards bring supercomputing to the masses.”
- Optimizing for parallelism improves the use of computing resources.
- And optimizing for data locality is often a first step!

The CPU-Memory Gap

The increasing gap between DRAM, disk, and CPU speeds.



Once upon a time, everything was slow in a computer.



The second space race ...

Who is the speaker?

Who are the attendees?

What are the prerequisites?

- Some familiarity with algorithms and their analysis.
- Elementary linear algebra (matrix multiplication).
- Ideas about multithreaded programming.
- Some ideas about GPUs.
- Visit <http://uwo.sharcnet.ca/>

What are the objectives of this tutorial?

- 1 Understand why data locality can have a huge impact on code performances.
- 2 Acquire some ideas on how data locality can be analyzed and improved.
- 3 Understand the concepts of work, span, parallelism, burdened parallelism in multithreaded programming.
- 4 Acquire some ideas on how parallelism can be analyzed and improved in multithreaded programming.
- 5 Understand issues related to memory bandwidth in GPU programming
- 6 Acquire some ideas on how the effective memory bandwidth of a GPU kernel can be improved.

Acknowledgments and references

Acknowledgments.

- Charles E. Leiserson (MIT), Matteo Frigo (Axis Semiconductor) Saman P. Amarasinghe (MIT) and Cyril Zeller (NVIDIA) for sharing with me the sources of their course notes and other documents.
- Yuzhen Xie (Maplesoft) and Anisul Sardar Haque (UWO) for their great help in the preparation of this tutorial.

References.

- *The Implementation of the Cilk-5 Multithreaded Language* by Matteo Frigo Charles E. Leiserson Keith H. Randall.
- *Cache-Oblivious Algorithms* by Matteo Frigo, Charles E. Leiserson, Harald Prokop and Sridhar Ramachandran.
- *The Cache Complexity of Multithreaded Cache Oblivious Algorithms* by Matteo Frigo and Volker Strumpfen.
- *How To Write Fast Numerical Code: A Small Introduction* by Srinivas Chellappa, Franz Franchetti, and Markus Pueschel.
- *Models of Computation: Exploring the Power of Computing* by John E. Savage.
- <http://developer.nvidia.com/category/zone/cuda-zone>
- <http://www.csd.uwo.ca/~moreno/HPC-Resources.html>

Plan

- 1 Data locality and cache misses
 - Hierarchical memories and their impact on our programs
 - Cache complexity and cache-oblivious algorithms put into practice
 - A detailed case study: counting sort
- 2 Multicore programming
 - Multicore architectures
 - Cilk / Cilk++ / Cilk Plus
 - The fork-join multithreaded programming model
- 3 GPU programming
 - The CUDA programming and memory models
 - Tiled matrix multiplication in CUDA
 - Optimizing Matrix Transpose with CUDA

Plan

- 1 Data locality and cache misses
 - Hierarchical memories and their impact on our programs
 - Cache complexity and cache-oblivious algorithms put into practice
 - A detailed case study: counting sort
- 2 Multicore programming
 - Multicore architectures
 - Cilk / Cilk++ / Cilk Plus
 - The fork-join multithreaded programming model
- 3 GPU programming
 - The CUDA programming and memory models
 - Tiled matrix multiplication in CUDA
 - Optimizing Matrix Transpose with CUDA

Capacity
Access Time
Cost

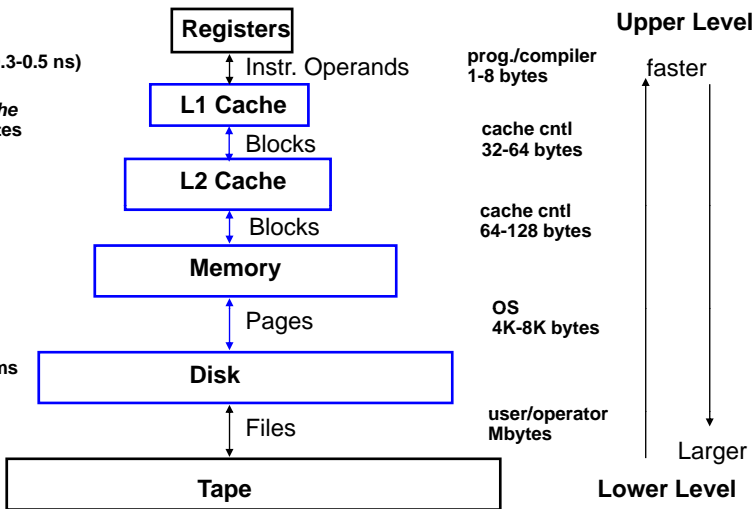
CPU Registers
 100s Bytes
 300 – 500 ps (0.3-0.5 ns)

L1 and L2 Cache
 10s-100s K Bytes
 ~1 ns - ~10 ns
 \$1000s/ GByte

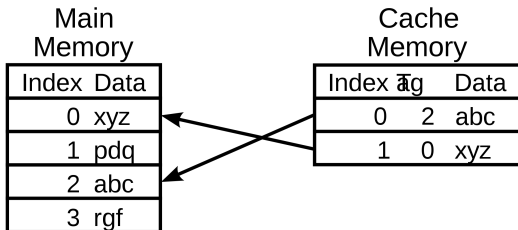
Main Memory
 G Bytes
 80ns- 200ns
 ~ \$100/ GByte

Disk
 10s T Bytes, 10 ms
 (10,000,000 ns)
 ~ \$1 / GByte

Tape
 infinite
 sec-min
 ~\$1 / GByte

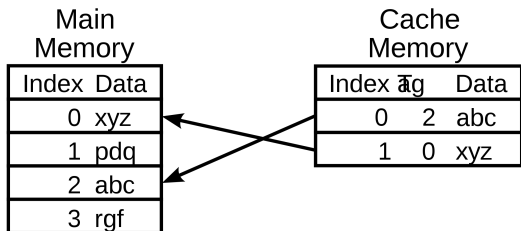


CPU Cache (1/7)



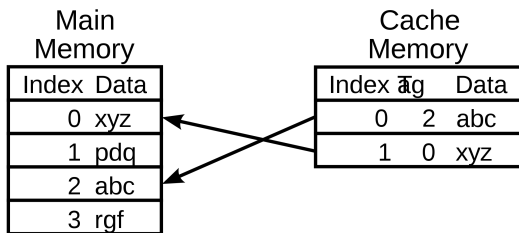
- A **CPU cache** is an auxiliary memory which is **smaller, faster memory** than the main memory and which stores **copies** of the main memory locations that are **expectedly frequently used**.
- Most modern desktop and server CPUs have at least three independent caches: the **data cache**, the **instruction cache** and the **translation look-aside buffer**.

CPU Cache (2/7)



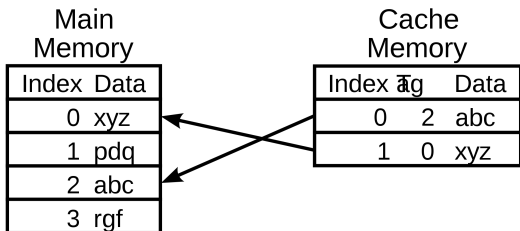
- Each location in each memory (main or cache) has
 - a datum (cache line) which ranges between 8 and 512 bytes in size, while a datum requested by a CPU instruction ranges between 1 and 16.
 - a unique index (called address in the case of the main memory)
- In the cache, each location has also a tag (storing the address of the corresponding cached datum).

CPU Cache (3/7)



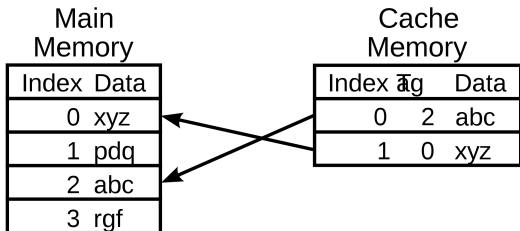
- When the CPU needs to read or write a location, it checks the cache:
 - if it finds it there, we have a **cache hit**
 - if not, we have a **cache miss** and (in most cases) the processor needs to create a new entry in the cache.
- Making room for a new entry requires a **replacement policy**: the **Least Recently Used** (LRU) discards the least recently used items first; this requires to use **age bits**.

CPU Cache (4/7)



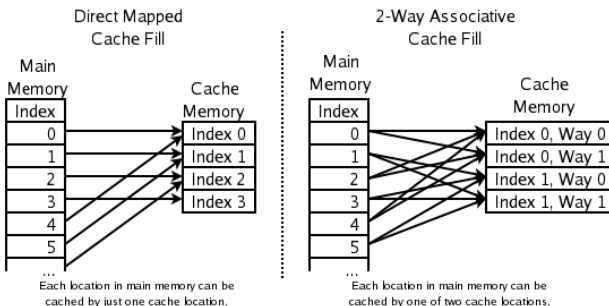
- **Read latency** (time to read a datum from the main memory) requires to keep the CPU busy with something else:
 - out-of-order execution:** attempt to execute independent instructions arising after the instruction that is waiting due to the cache miss
 - hyper-threading (HT):** allows an alternate thread to use the CPU

CPU Cache (5/7)

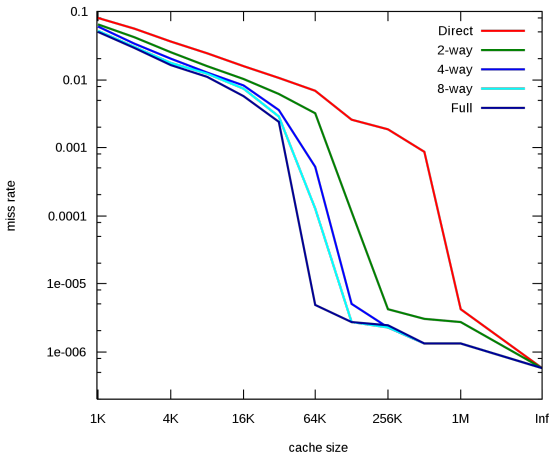


- Modifying data in the cache requires a **write policy** for updating the main memory
 - **write-through cache**: writes are immediately mirrored to main memory
 - **write-back cache**: the main memory is mirrored when that data is evicted from the cache
- The cache copy may become out-of-date or stale, if other processors modify the original entry in the main memory.

CPU Cache (6/7)



- The replacement policy decides where in the cache a copy of a particular entry of main memory will go:
 - **fully associative**: any entry in the cache can hold it
 - **direct mapped**: only one possible entry in the cache can hold it
 - **N -way set associative**: N possible entries can hold it



- *Cache Performance for SPEC CPU2000* by J.F. Cantin and M.D. Hill.
- The SPEC CPU2000 suite is a collection of 26 compute-intensive, non-trivial programs used to evaluate the performance of a computer's CPU, memory system, and compilers (<http://www.spec.org/osg/cpu2000>).

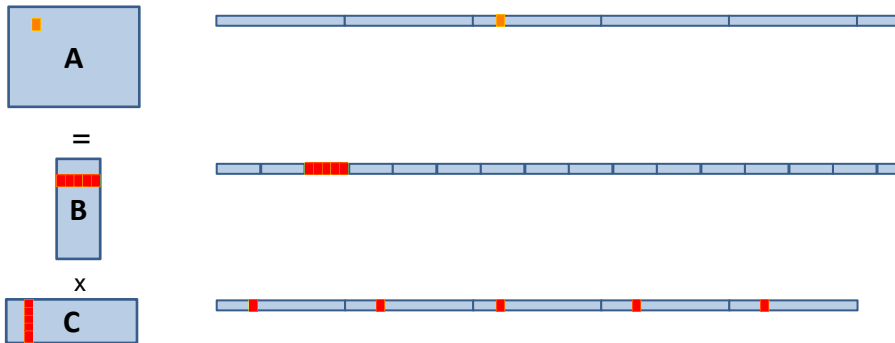
Cache issues

- **Cold miss:** The first time the data is available. Cure: Prefetching may be able to reduce this type of cost.
- **Capacity miss:** The previous access has been evicted because too much data touched in between, since the *working data set* is too large. Cure: Reorganize the data access such that *reuse* occurs before eviction.
- **Conflict miss:** Multiple data items mapped to the same location with eviction before cache is full. Cure: Rearrange data and/or pad arrays.
- **True sharing miss:** Occurs when a thread in another processor wants the same data. Cure: Minimize sharing.
- **False sharing miss:** Occurs when another processor uses different data in the same cache line. Cure: Pad data.

A typical matrix multiplication C code

```
#define IND(A, x, y, d) A[(x)*(d)+(y)]
uint64_t testMM(const int x, const int y, const int z)
{
    double *A; double *B; double *C;
    long started, ended;
    float timeTaken;
    int i, j, k;
    srand(getSeed());
    A = (double *)malloc(sizeof(double)*x*y);
    B = (double *)malloc(sizeof(double)*x*z);
    C = (double *)malloc(sizeof(double)*y*z);
    for (i = 0; i < x*z; i++) B[i] = (double) rand() ;
    for (i = 0; i < y*z; i++) C[i] = (double) rand() ;
    for (i = 0; i < x*y; i++) A[i] = 0 ;
    started = example_get_time();
    for (i = 0; i < x; i++)
        for (j = 0; j < y; j++)
            for (k = 0; k < z; k++)
                // A[i][j] += B[i][k] + C[k][j];
                IND(A,i,j,y) += IND(B,i,k,z) * IND(C,k,j,z);
    ended = example_get_time();
    timeTaken = (ended - started)/1.f;
    return timeTaken;
}
```

Issues with matrix representation

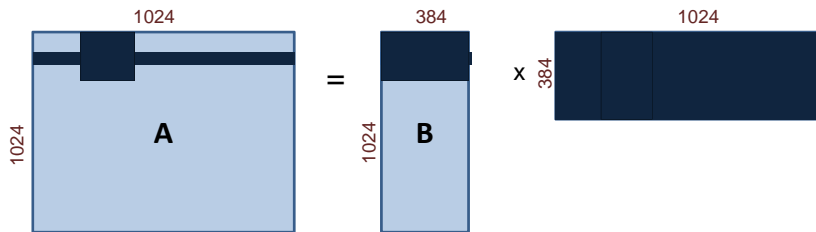


- Contiguous accesses are better:
 - Data fetch as cache line (Core 2 Duo 64 byte per cache line)
 - With contiguous data, a single cache fetch supports 8 reads of doubles.
 - **Transposing the matrix C should reduce L1 cache misses!**

Transposing for optimizing spatial locality

```
float testMM(const int x, const int y, const int z)
{
    double *A; double *B; double *C; double *Cx;
    long started, ended; float timeTaken; int i, j, k;
    A = (double *)malloc(sizeof(double)*x*y);
    B = (double *)malloc(sizeof(double)*x*z);
    C = (double *)malloc(sizeof(double)*y*z);
    Cx = (double *)malloc(sizeof(double)*y*z);
    srand(getSeed());
    for (i = 0; i < x*z; i++) B[i] = (double) rand() ;
    for (i = 0; i < y*z; i++) C[i] = (double) rand() ;
    for (i = 0; i < x*y; i++) A[i] = 0 ;
    started = example_get_time();
    for(j =0; j < y; j++)
        for(k=0; k < z; k++)
            IND(Cx,j,k,z) = IND(C,k,j,y);
    for (i = 0; i < x; i++)
        for (j = 0; j < y; j++)
            for (k = 0; k < z; k++)
                IND(A, i, j, y) += IND(B, i, k, z) *IND(Cx, j, k, z);
    ended = example_get_time();
    timeTaken = (ended - started)/1.f;
    return timeTaken;
}
```

Issues with data reuse



- Naive calculation of a row of A, so computing 1024 coefficients: 1024 accesses in A, 384 in B and $1024 \times 384 = 393,216$ in C. Total = 394,524.
- Computing a 32×32 -block of A, so computing again 1024 coefficients: 1024 accesses in A, 384×32 in B and 32×384 in C. Total = 25,600.
- The iteration space is traversed so as to reduce memory accesses.

Blocking for optimizing temporal locality

```
float testMM(const int x, const int y, const int z)
{
    double *A; double *B; double *C;
    long started, ended; float timeTaken; int i, j, k, i0, j0, k0;
    A = (double *)malloc(sizeof(double)*x*y);
    B = (double *)malloc(sizeof(double)*x*z);
    C = (double *)malloc(sizeof(double)*y*z);
    srand(getSeed());
    for (i = 0; i < x*z; i++) B[i] = (double) rand() ;
    for (i = 0; i < y*z; i++) C[i] = (double) rand() ;
    for (i = 0; i < x*y; i++) A[i] = 0 ;
    started = example_get_time();
    for (i = 0; i < x; i += BLOCK_X)
        for (j = 0; j < y; j += BLOCK_Y)
            for (k = 0; k < z; k += BLOCK_Z)
                for (i0 = i; i0 < min(i + BLOCK_X, x); i0++)
                    for (j0 = j; j0 < min(j + BLOCK_Y, y); j0++)
                        for (k0 = k; k0 < min(k + BLOCK_Z, z); k0++)
                            IND(A,i0,j0,y) += IND(B,i0,k0,z) * IND(C,k0,j0,y);
    ended = example_get_time();
    timeTaken = (ended - started)/1.f;
    return timeTaken;
}
```

Transposing and blocking for optimizing data locality

```
float testMM(const int x, const int y, const int z)
{
    double *A; double *B; double *C;
    long started, ended; float timeTaken; int i, j, k, i0, j0, k0;
    A = (double *)malloc(sizeof(double)*x*y);
    B = (double *)malloc(sizeof(double)*x*z);
    C = (double *)malloc(sizeof(double)*y*z);
    srand(getSeed());
    for (i = 0; i < x*z; i++) B[i] = (double) rand() ;
    for (i = 0; i < y*z; i++) C[i] = (double) rand() ;
    for (i = 0; i < x*y; i++) A[i] = 0 ;
    started = example_get_time();
    for (i = 0; i < x; i += BLOCK_X)
        for (j = 0; j < y; j += BLOCK_Y)
            for (k = 0; k < z; k += BLOCK_Z)
                for (i0 = i; i0 < min(i + BLOCK_X, x); i0++)
                    for (j0 = j; j0 < min(j + BLOCK_Y, y); j0++)
                        for (k0 = k; k0 < min(k + BLOCK_Z, z); k0++)
                            IND(A,i0,j0,y) += IND(B,i0,k0,z) * IND(C,j0,k0,z);
    ended = example_get_time();
    timeTaken = (ended - started)/1.f;
    return timeTaken;
}
```

Experimental results

Computing the product of two $n \times n$ matrices on my laptop (Core2 Duo CPU P8600 @ 2.40GHz, L1 cache of 3072 KB, 4 GBytes of RAM)

n	naive	transposed	speedup	64×64 -tiled	speedup	t. & t.	speedup
128	7	3		7		2	
256	26	43		155		23	
512	1805	265	6.81	1928	0.936	187	9.65
1024	24723	3730	6.62	14020	1.76	1490	16.59
2048	271446	29767	9.11	112298	2.41	11960	22.69
4096	2344594	238453	9.83	1009445	2.32	101264	23.15

Timings are in milliseconds.

The cache-oblivious multiplication (more on this later) runs within 12978 and 106758 for $n = 2048$ and $n = 4096$ respectively.

Other performance counters

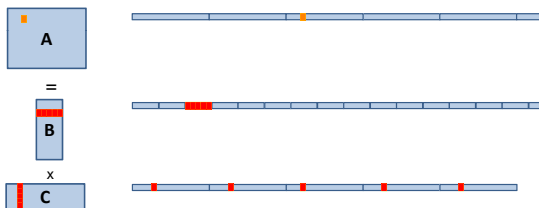
Hardware count events

- **CPI Clock cycles Per Instruction:** the number of clock cycles that happen when an instruction is being executed. With pipelining we can improve the CPI by exploiting instruction level parallelism
- **L1 and L2 Cache Miss Rate.**
- **Instructions Retired:** In the event of a misprediction, instructions that were scheduled to execute along the mispredicted path must be canceled.

	CPI	L1 Miss Rate	L2 Miss Rate	Percent SSE Instructions	Instructions Retired
In C	4.78	0.24	0.02	43%	13,137,280,000
Transposed	1.13	0.15	0.02	50%	13,001,486,336
Tiled	0.49	0.02	0	39%	18,044,811,264

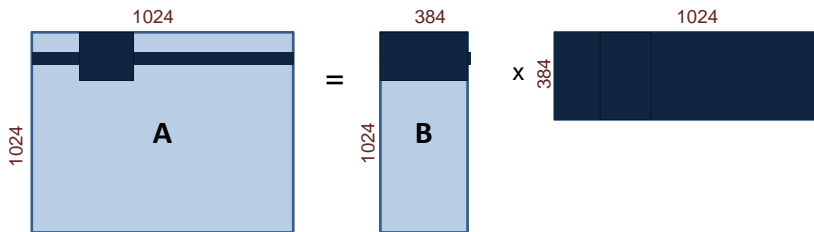
Annotations from image:
 - CPI: In C (4.78) is 5x Transposed (1.13) and 3x Tiled (0.49).
 - L1 Miss Rate: In C (0.24) is 2x Transposed (0.15) and 8x Tiled (0.02).
 - Instructions Retired: In C (13,137,280,000) is 1x Transposed (13,001,486,336) and 0.8x Tiled (18,044,811,264).

Analyzing cache misses in the naive and transposed multiplication



- Let A , B and C have format (m, n) , (m, p) and (p, n) respectively.
- A is scanned once, so mn/L cache misses if L is the number of coefficients per cache line.
- B is scanned n times, so mnp/L cache misses if the cache cannot hold a row.
- C is accessed “nearly randomly” (for m large enough) leading to mnp cache misses.
- Since $2mnp$ arithmetic operations are performed, this means roughly **one cache miss per flop!**
- If C is transposed, then the ratio improves to 1 for L .

Analyzing cache misses in the tiled multiplication



- Let A , B and C have format (m, n) , (m, p) and (p, n) respectively.
- Assume all tiles are square of order b and three fit in cache.
- If C is transposed, then loading three blocks in cache cost $3b^2/L$.
- This process happens n^3/b^3 times, leading to $3n^3/(bL)$ cache misses.
- Three blocks fit in cache for $3b^2 < Z$, if Z is the cache size.
- So $O(n^3/(\sqrt{Z}L))$ cache misses, if b is well chosen, which is optimal.

Plan

- 1 Data locality and cache misses
 - Hierarchical memories and their impact on our programs
 - Cache complexity and cache-oblivious algorithms put into practice
 - A detailed case study: counting sort
- 2 Multicore programming
 - Multicore architectures
 - Cilk / Cilk++ / Cilk Plus
 - The fork-join multithreaded programming model
- 3 GPU programming
 - The CUDA programming and memory models
 - Tiled matrix multiplication in CUDA
 - Optimizing Matrix Transpose with CUDA

The (Z, L) ideal cache model (1/4)

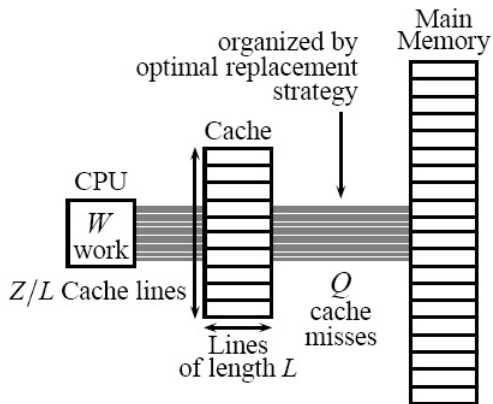


Figure 1: The ideal-cache model

The (Z, L) ideal cache model (2/4)

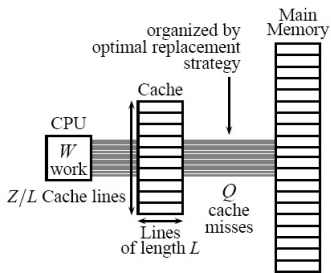


Figure 1: The ideal-cache model

- Computer with a **two-level memory hierarchy**:
 - an ideal (data) cache of Z words partitioned into Z/L cache lines, where L is the number of words per cache line.
 - an arbitrarily large main memory.
- Data moved between cache and main memory are always cache lines.
- The cache is **tall**, that is, Z is much larger than L , say $Z \in \Omega(L^2)$.

The (Z, L) ideal cache model (3/4)

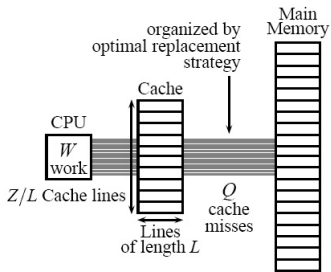


Figure 1: The ideal-cache model

- The processor can only reference words that reside in the cache.
- If the referenced word belongs to a line already in cache, a **cache hit** occurs, and the word is delivered to the processor.
- Otherwise, a **cache miss** occurs, and the line is fetched into the cache.

The (Z, L) ideal cache model (4/4)

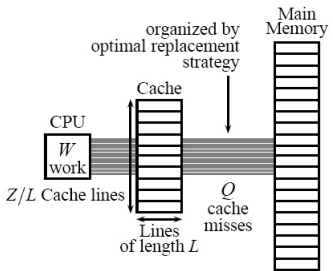


Figure 1: The ideal-cache model

- The ideal cache is **fully associative**: cache lines can be stored anywhere in the cache.
- The ideal cache uses the **optimal off-line strategy of replacing** the cache line whose next access is furthest in the future, and thus it exploits temporal locality perfectly.

Cache complexity

- For an algorithm with an input of size n , the ideal-cache model uses two complexity measures:
 - the **work complexity** $W(n)$, which is its conventional running time in a RAM model.
 - the **cache complexity** $Q(n; Z, L)$, the number of cache misses it incurs (as a function of the size Z and line length L of the ideal cache).
 - When Z and L are clear from context, we simply write $Q(n)$ instead of $Q(n; Z, L)$.
- An algorithm is said to be **cache aware** if its behavior (and thus performances) can be tuned (and thus depend on) on the particular cache size and line length of the targeted machine.
- Otherwise the algorithm is **cache oblivious**.

Cache complexity of an array scanning

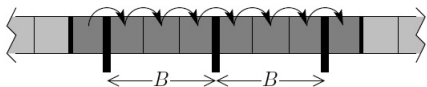


Figure 2. Scanning an array of N elements arbitrarily aligned with blocks may cost one more memory transfer than $\lceil N/B \rceil$.

- B and N on the picture are our L and n .
- Consider an array of n words in main memory.
- Loading its elements by **scanning** incurs $n/L + 1$ cache misses.
- That would be n/L if L divides L and the array is aligned, that is, starts and ends with a cache line.
- We will often use this remark and, for simplicity, we will always neglect the $+1$.

Cache complexity of the naive and tiled matrix multiplications

- Consider square matrices of order n and an (Z, L) -ideal cache.
- The naive multiplication (as specified before)

```

for(i =0; i < n; i++)
    for(j =0; j < n; j++)
        for(k=0; k < n; k++)
            C[i][j] += A[i][k] * B[k][j];

```

incurs $O(n^3)$ cache misses, for n large enough ($n > Z^2$).

- The tiled multiplication (as specified before)

```

for(i =0; i < n/s; i++)
    for(j =0; j < n/s; j++)
        for(k=0; k < n/s; k++)
            blockMult(A,B,C,i,j,k,s);

```

incurs $\Theta(n^3/(L\sqrt{Z}))$ cache misses, for n large enough ($n > \sqrt{Z}$) which can be proved to be optimal, though cache-aware.

A matrix transposition cache-oblivious and cache-optimal algorithm

- Given an $m \times n$ matrix A stored in a row-major layout, compute and store A^T into an $n \times m$ matrix B also stored in a row-major layout.
- A naive approach would incur $O(mn)$ cache misses, for n, m large enough.
- The algorithm REC-TRANSPOSE below incurs $\Theta(1 + mn/L)$ cache misses, which is optimal.
 - If $n \geq m$, the REC-TRANSPOSE algorithm partitions

$$A = (A_1 \ A_2) , \quad B = \begin{pmatrix} B_1 \\ B_2 \end{pmatrix}$$

and recursively executes REC-TRANSPOSE(A_1, B_1) and REC-TRANSPOSE(A_2, B_2).

- If $m > n$, the REC-TRANSPOSE algorithm partitions

$$A = \begin{pmatrix} A_1 \\ A_2 \end{pmatrix} , \quad B = (B_1 \ B_2)$$

and recursively executes REC-TRANSPOSE(A_1, B_1) and REC-TRANSPOSE(A_2, B_2).

Cache-oblivious matrix transposition works in practice!

size	Naive	Cache-oblivious	ratio
5000x5000	126	79	1.59
10000x10000	627	311	2.02
20000x20000	4373	1244	3.52
30000x30000	23603	2734	8.63
40000x40000	62432	4963	12.58

- Intel(R) Xeon(R) CPU E7340 @ 2.40GHz
- L1 data 32 KB, L2 4096 KB, cache line size 64bytes
- **Both codes run on 1 core** on a node with 128GB.
- The ration comes simply from an **optimal memory access pattern**.

A cache-oblivious matrix multiplication algorithm

- To multiply an $m \times n$ matrix A and an $n \times p$ matrix B , the REC-MULT algorithm halves the largest of the three dimensions and recurs according to one of the following three cases:

$$\begin{pmatrix} A_1 \\ A_2 \end{pmatrix} B = \begin{pmatrix} A_1 B \\ A_2 B \end{pmatrix}, \quad (1)$$

$$(A_1 \ A_2) \begin{pmatrix} B_1 \\ B_2 \end{pmatrix} = A_1 B_1 + A_2 B_2, \quad (2)$$

$$A (B_1 \ B_2) = (AB_1 \ AB_2). \quad (3)$$

- In case (1), we have $m \geq \max\{n, p\}$. Matrix A is split horizontally, and both halves are multiplied by matrix B .
- In case (2), we have $n \geq \max\{m, p\}$. Both matrices are split, and the two halves are multiplied.
- In case (3), we have $p \geq \max\{m, n\}$. Matrix B is split vertically, and each half is multiplied by A .
- The base case occurs when $m = n = p = 1$.
- The algorithm REC-MULT above incurs $\Theta(m + n + p + (mn + np + mp)/L + mnp/(L\sqrt{Z}))$ cache misses, which is optimal.

Plan

- 1 Data locality and cache misses
 - Hierarchical memories and their impact on our programs
 - Cache complexity and cache-oblivious algorithms put into practice
 - A detailed case study: counting sort
- 2 Multicore programming
 - Multicore architectures
 - Cilk / Cilk++ / Cilk Plus
 - The fork-join multithreaded programming model
- 3 GPU programming
 - The CUDA programming and memory models
 - Tiled matrix multiplication in CUDA
 - Optimizing Matrix Transpose with CUDA

Counting sort: the algorithm

- *Counting sort* takes as input a collection of n items, each of which known by a key in the range $0 \dots k$.
- The algorithm computes a *histogram* of the number of times each key occurs.
- Then performs a *prefix sum* to compute positions in the output.

```
allocate an array Count[0..k]; initialize each array cell to zero
for each input item x:
    Count[key(x)] = Count[key(x)] + 1
total = 0
for i = 0, 1, ... k:
    c = Count[i]
    Count[i] = total
    total = total + c
allocate an output array Output[0..n-1]
for each input item x:
    store x in Output[Count[key(x)]]
    Count[key(x)] = Count[key(x)] + 1
return Output
```

Counting sort: cache complexity analysis

```
allocate an array Count[0..k]; initialize each array cell to zero
for each input item x:
    Count[key(x)] = Count[key(x)] + 1
total = 0
for i = 0, 1, ... k:
    c = Count[i]
    Count[i] = total
    total = total + c
allocate an output array Output[0..n-1]
for each input item x:
    store x in Output[Count[key(x)]]
    Count[key(x)] = Count[key(x)] + 1
return Output
```

- 1 n/L to compute k .
 - 2 k/L cache misses to initialize Count.
 - 3 $n/L + n$ cache misses for the histogram (worst case).
 - 4 k/L cache misses for the prefix sum.
 - 5 $n/L + n + n$ cache misses for building Output (worst case).
- Total: $3n + 3n/L + 2k/L$ cache misses (worst case).

Counting sort: cache complexity analysis: explanations

- 1 n/L to compute k : this can be done by traversing the `items` linearly.
 - 2 k/L cache misses to initialize `Count`: this can be done by traversing the `Count` linearly.
 - 3 $n/L + n$ cache misses for the histogram (worst case): `items` accesses are linear but `Count` accesses are potentially random.
 - 4 k/L cache misses for the prefix sum: `Count` accesses are linear.
 - 5 $n/L + n + n$ cache misses for building `Output` (worst case): `items` accesses are linear but `Output` and `Count` accesses are potentially random.
- Total: $3n + 3n/L + 2k/L$ cache misses (worst case).

How to fix the poor data locality of counting sort?

```
allocate an array Count[0..k]; initialize each array cell to zero
for each input item x:
    Count[key(x)] = Count[key(x)] + 1
total = 0
for i = 0, 1, ... k:
    c = Count[i]
    Count[i] = total
    total = total + c
allocate an output array Output[0..n-1]
for each input item x:
    store x in Output[Count[key(x)]]
    Count[key(x)] = Count[key(x)] + 1
return Output
```

- Recall that our worst case is $3n + 3n/L + 2k/L$ cache misses.
- The troubles come from the irregular memory accesses which experience **capacity misses** and **conflict misses**.
- Workaround: we preprocess the input so that counting sort is applied in succession to several smaller input sets with smaller value ranges.
- To put it simply, so that k and n are small enough for Output and Count to incur cold misses only.

Counting sort: bucketing the input

```

allocate an array bucketsize[0..m-1]; initialize each array cell to zero
for each input item x:
    bucketsize[floor(key(x) m/(k+1))] := bucketsize[floor(key(x) m/(k+1))] + 1
total = 0
for i = 0, 1, ... m-1:
    c = bucketsize[i]
    bucketsize[i] = total
    total = total + c
allocate an array bucketedinput[0..n-1];
for each input item x:
    q := floor(key(x) m/(k+1))
    bucketedinput[bucketsize[q] ] := key(x)
    bucketsize[q] := bucketsize[q] + 1
return bucketedinput

```

- Goal: after preprocessing, Count and Output incur **cold misses only**.
- To this end we choose a parameter m (more on this later) such that
 - 1 a key in the range $[ih, (i+1)h - 1]$ is always before a key in the range $[(i+1)h, (i+2)h - 1]$, for $i = 0 \dots m-2$, with $h = k/m$,
 - 2 bucketsize and m cache-lines from bucketedinput all fit in cache. That is, counting cache-lines, $m/L + m \leq Z/L$, that is, $m + mL \leq Z$.

Counting sort: cache complexity with bucketing

```

allocate an array bucketsize[0..m-1]; initialize each array cell to zero
for each input item x:
    bucketsize[floor(key(x) m/(k+1))] := bucketsize[floor(key(x) m/(k+1))] + 1
total = 0
for i = 0, 1, ... m-1:
    c = bucketsize[i]
    bucketsize[i] = total
    total = total + c
allocate an array bucketedinput[0..n-1];
for each input item x:
    q := floor(key(x) m/(k+1))
    bucketedinput[bucketsize[q] ] := key(x)
    bucketsize[q] := bucketsize[q] + 1
return bucketedinput

```

- ① $3m/L + n/L$ cache misses to compute bucketsize
 - ② **Key observation:** bucketedinput is traversed regularly by segment.
 - ③ Hence, $2n/L + m + m/L$ cache misses to compute bucketedinput
- Preprocessing: $3n/L + 3m/L + m$ cache misses.

Counting sort: cache complexity with bucketing: explanations

- ① $3m/L + n/L$ caches misses to compute bucketsize:
 - m/L to set each cell of bucketsize to zero,
 - $m/L + n/L$ for the first for loop,
 - m/L for the second for loop.
- ② **Key observation:** bucketedinput is traversed regularly by segment:
 - So writing bucketedinput means writing (in a linear traversal) m consecutive arrays, of possibly different sizes, but with total size n .
 - Thus, because of possible misalignments between those arrays and their cache-lines, this writing procedure can yield $n/L + m$ cache misses (and not just n/L).
- ③ Hence, $2n/L + m + m/L$ caches misses to compute bucketedinput:
 - n/L to read the items,
 - $n/L + m$ to write bucketedinput,
 - m/L to load bucketsize.

Cache friendly counting sort: complete cache complexity analysis

- **Assumption:** the preprocessing creates buckets of average size n/m .
- After preprocessing, counting sort is applied to each bucket whose values are in a range $[ih, (i+1)h - 1]$, for $i = 0 \dots m - 1$.
- To be cache-friendly, this requires, for $i = 0 \dots m - 1$, $h + |\{\text{key} \in [ih, (i+1)h - 1]\}| < Z$ and $m < Z/(1+L)$. These two are very realistic assumption considering today's cache size.
- And the total complexity becomes;

$$\begin{aligned}
 Q_{\text{total}} &= Q_{\text{preprocessing}} + Q_{\text{sorting}} \\
 &= Q_{\text{preprocessing}} + m Q_{\text{sorting of one bucket}} \\
 &= Q_{\text{preprocessing}} + m \left(3 \frac{n}{mL} + 3 \frac{n}{m} + 2 \frac{k}{mL} \right) \\
 &= Q_{\text{preprocessing}} + 6n/L + 2k/L \\
 &= 3n/L + 3m/L + m + 6n/L + 2k/L \\
 &= 9n/L + 3m/L + m + 2k/L
 \end{aligned}$$

Instead of $3n + 3n/L + 2k/L$ for the naive counting sort.

Cache friendly counting sort: experimental results

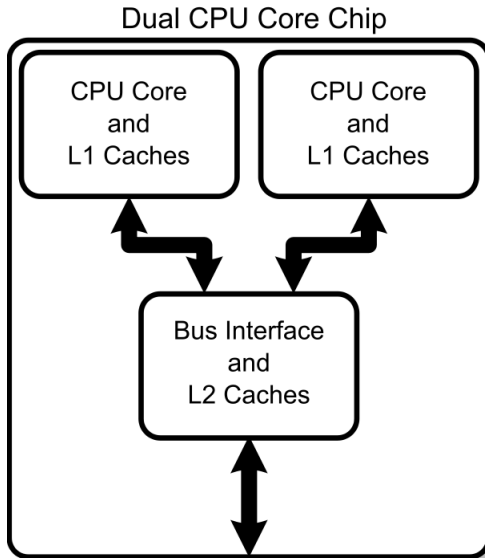
- Experimentation on an *Intel(R) Core(TM) i7 CPU @ 2.93GHz*. It has L2 cache of 8MB.
- CPU times in seconds for both classical and cache-friendly counting sort algorithm.
- The keys are random machine integers in the range $[0, n]$.

n	classical counting sort	cache-oblivious counting sort (preprocessing + sorting)
100000000	13.74	4.66 (3.04 + 1.62)
200000000	30.20	9.93 (6.16 + 3.77)
300000000	50.19	16.02 (9.32 + 6.70)
400000000	71.55	22.13 (12.50 + 9.63)
500000000	94.32	28.37 (15.71 + 12.66)
600000000	116.74	34.61 (18.95 + 15.66)

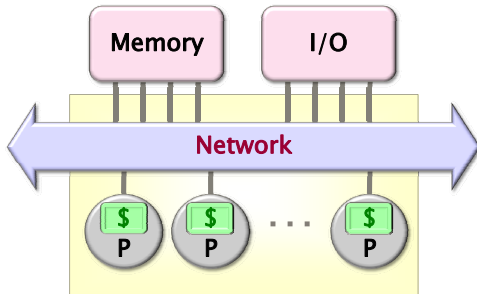
Summary and notes

Plan

- 1 Data locality and cache misses
 - Hierarchical memories and their impact on our programs
 - Cache complexity and cache-oblivious algorithms put into practice
 - A detailed case study: counting sort
- 2 Multicore programming
 - Multicore architectures
 - Cilk / Cilk++ / Cilk Plus
 - The fork-join multithreaded programming model
- 3 GPU programming
 - The CUDA programming and memory models
 - Tiled matrix multiplication in CUDA
 - Optimizing Matrix Transpose with CUDA



- A **multi-core processor** is an integrated circuit to which two or more individual processors (called cores in this sense) have been attached.



Chip Multiprocessor (CMP)

- Cores on a multi-core device can be **coupled tightly or loosely**:
 - may share or may not share a cache,
 - implement inter-core communications methods or message passing.
- Cores on a multi-core implement the **same architecture features as single-core systems** such as instruction pipeline parallelism (ILP), vector-processing, hyper-threading, etc.

Cache Coherence (1/6)

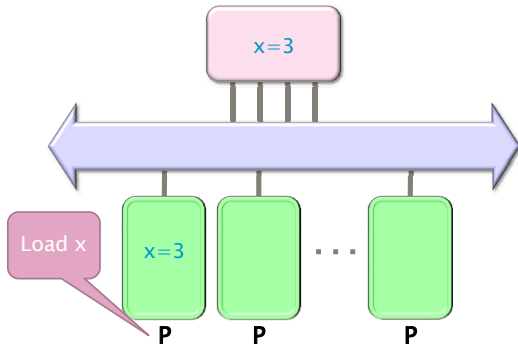


Figure: Processor P_1 reads $x=3$ first from the backing store (higher-level memory)

Cache Coherence (2/6)

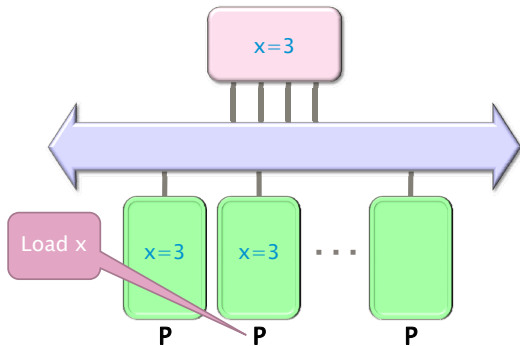


Figure: Next, Processor P_2 loads $x=3$ from the same memory

Cache Coherence (3/6)

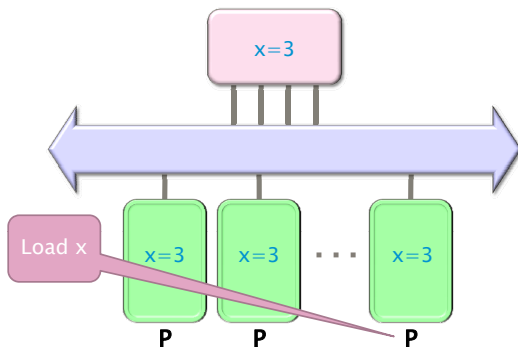


Figure: Processor P_4 loads $x=3$ from the same memory

Cache Coherence (4/6)

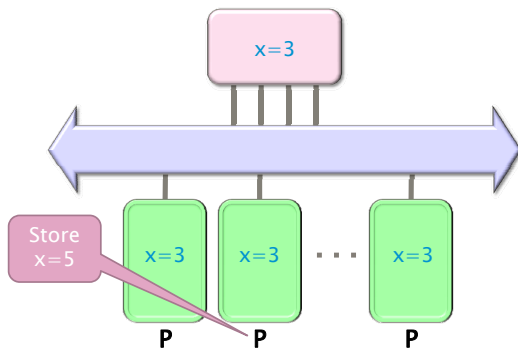


Figure: Processor P_2 issues a write $x=5$

Cache Coherence (5/6)

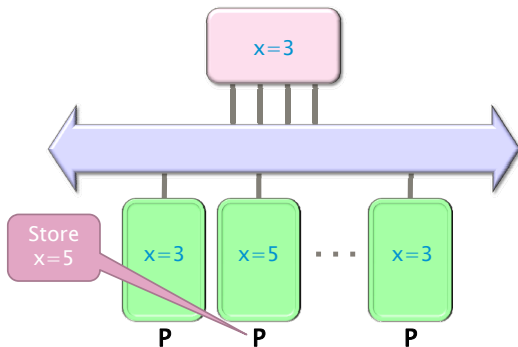


Figure: Processor P_2 writes $x=5$ in his local cache

Cache Coherence (6/6)

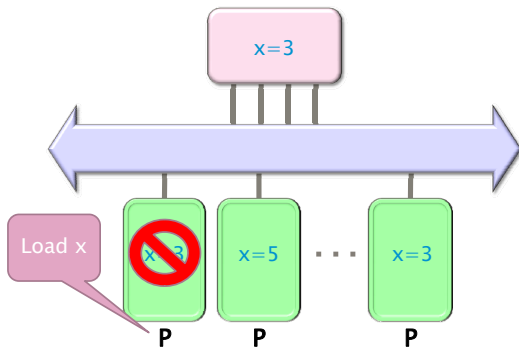


Figure: Processor P_1 issues a read x , which is now invalid in its cache

MSI Protocol

- In this cache coherence protocol each block contained inside a cache can have one of three possible states:
 - **M**: the cache line has been **modified** and the corresponding data is inconsistent with the backing store; the cache has the responsibility to write the block to the backing store when it is evicted.
 - **S**: this block is unmodified and is **shared**, that is, exists in at least one cache. The cache can evict the data without writing it to the backing store.
 - **I**: this block is **invalid**, and must be fetched from memory or another cache if the block is to be stored in this cache.
- These coherency states are maintained through communication between the caches and the backing store.
- The caches have different responsibilities when blocks are read or written, or when they learn of other caches issuing reads or writes for a block.

True Sharing and False Sharing

- **True sharing:**

- True sharing cache misses occur whenever two processors access the same data word
- True sharing requires the processors involved to explicitly synchronize with each other to ensure program correctness.
- A computation is said to have **temporal locality** if it re-uses much of the data it has been accessing.
- Programs with high temporal locality tend to have less true sharing.

- **False sharing:**

- False sharing results when different processors use different data that happen to be co-located on the same cache line
 - A computation is said to have **spatial locality** if it uses multiple words in a cache line before the line is displaced from the cache
 - Enhancing spatial locality often minimizes false sharing
- See *Data and Computation Transformations for Multiprocessors* by J.M. Anderson, S.P. Amarasinghe and M.S. Lam
<http://suif.stanford.edu/papers/anderson95/paper.html>

Multi-core processor (cntd)

- **Advantages:**

- Cache coherency circuitry operate at higher rate than off-chip.
- Reduced power consumption for a dual core vs two coupled single-core processors (better quality communication signals, cache can be shared)

- **Challenges:**

- Adjustments to existing software (including OS) are required to maximize performance
- Production yields down (an Intel quad-core is in fact a double dual-core)
- Two processing cores sharing the same bus and memory bandwidth may limit performances
- High levels of false or true sharing and synchronization can easily overwhelm the advantage of parallelism

Plan

- 1 Data locality and cache misses
 - Hierarchical memories and their impact on our programs
 - Cache complexity and cache-oblivious algorithms put into practice
 - A detailed case study: counting sort
- 2 Multicore programming
 - Multicore architectures
 - Cilk / Cilk++ / Cilk Plus
 - The fork-join multithreaded programming model
- 3 GPU programming
 - The CUDA programming and memory models
 - Tiled matrix multiplication in CUDA
 - Optimizing Matrix Transpose with CUDA

From Cilk to Cilk++ and Cilk Plus

- Cilk has been developed since 1994 at the MIT Laboratory for Computer Science by Prof. Charles E. Leiserson and his group, in particular by Matteo Frigo.
- Besides being used for research and teaching, Cilk was the system used to code the three world-class chess programs: Tech, Socrates, and Cilkchess.
- Over the years, the implementations of Cilk have run on computers ranging from networks of Linux laptops to an 1824-nodes Intel Paragon.
- From 2007 to 2009 Cilk has led to Cilk++, developed by Cilk Arts, an MIT spin-off, which was acquired by Intel in July 2009 and became Cilk Plus, see <http://www.cilk.com/>
- Cilk++ can be freely downloaded at <http://software.intel.com/en-us/articles/download-intel-cilk>
- Cilk is still developed at MIT <http://supertech.csail.mit.edu/cilk/>

Cilk++ (and Cilk Plus)

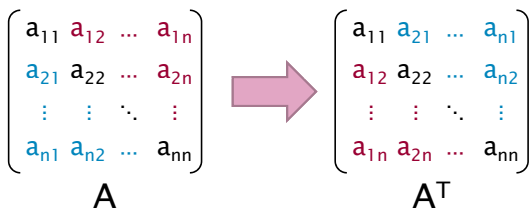
- Cilk++ (resp. Cilk) is a **small set of linguistic extensions to C++** (resp. C) supporting **fork-join parallelism**
- Both Cilk and Cilk++ feature a **provably efficient work-stealing scheduler**.
- Cilk++ provides a **hyperobject library** for parallelizing code with global variables and performing reduction for data aggregation.
- Cilk++ includes the **Cilkscreen** race detector and the **Cilkview** performance analyzer.

Nested Parallelism in Cilk ++

```
int fib(int n)
{
    if (n < 2) return n;
    int x, y;
    x = cilk_spawn fib(n-1);
    y = fib(n-2);
    cilk_sync;
    return x+y;
}
```

- The named **child** function `cilk_spawn fib(n-1)` may execute in parallel with its **parent**
- Cilk++ keywords `cilk_spawn` and `cilk_sync` grant **permissions for parallel execution**. They do not command parallel execution.

Loop Parallelism in Cilk ++



```
// indices run from 0, not 1
cilk_for (int i=1; i<n; ++i) {
    for (int j=0; j<i; ++j) {
        double temp = A[i][j];
        A[i][j] = A[j][i];
        A[j][i] = temp;
    }
}
```

The iterations of a `cilk_for` loop may execute in parallel.

Serial Semantics (1/2)

- Cilk (resp. Cilk++) is a multithreaded language for parallel programming that generalizes the semantics of C (resp. C++) by introducing linguistic constructs for parallel control.
- Cilk (resp. Cilk++) is a **faithful extension** of C (resp. C++):
 - The C (resp. C++) elision of a Cilk (resp. Cilk++) is a correct implementation of the semantics of the program.
 - Moreover, on one processor, a parallel Cilk (resp. Cilk++) program scales down to run nearly as fast as its C (resp. C++) elision.
- To obtain the serialization of a Cilk++ program

```
#define cilk_for for
#define cilk_spawn
#define cilk_sync
```


Serial Semantics (2/2)

```
int fib (int n) {  
  if (n<2) return (n);  
  else {  
    int x,y;  
    x = cilk_spawn fib(n-1);  
    y = fib(n-2);  
    cilk_sync;  
    return (x+y);  
  }  
}
```

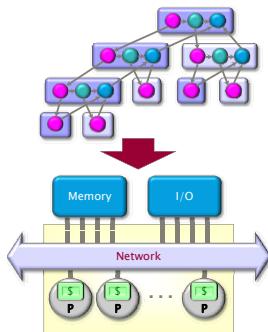
Cilk++ source

↓

```
int fib (int n) {  
  if (n<2) return (n);  
  else {  
    int x,y;  
    x = fib(n-1);  
    y = fib(n-2);  
    return (x+y);  
  }  
}
```

Serialization

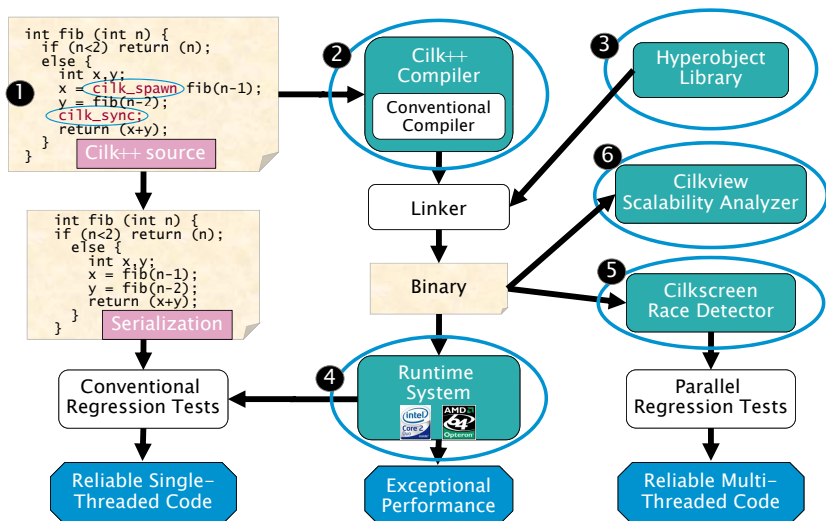
Scheduling



A **scheduler**'s job is to map a computation to particular processors. Such a mapping is called a **schedule**.

- If decisions are made at runtime, the scheduler is *online*, otherwise, it is *offline*
- Cilk++'s scheduler maps strands onto processors dynamically at runtime.

The Cilk++ Platform



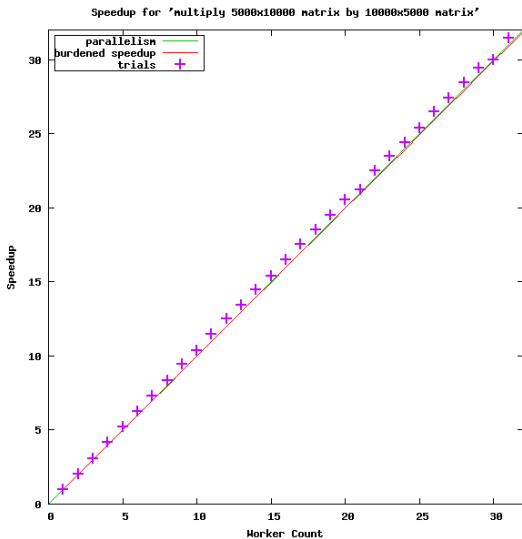
Benchmarks for the parallel version of the cache-oblivious mm

Multiplying a 4000x8000 matrix by a 8000x4000 matrix

- on 32 cores = 8 sockets x 4 cores (Quad Core AMD Opteron 8354) per socket.
- The 32 cores share a L3 32-way set-associative cache of 2 Mbytes.

#core	Elision (s)	Parallel (s)	speedup
8	420.906	51.365	8.19
16	432.419	25.845	16.73
24	413.681	17.361	23.83
32	389.300	13.051	29.83

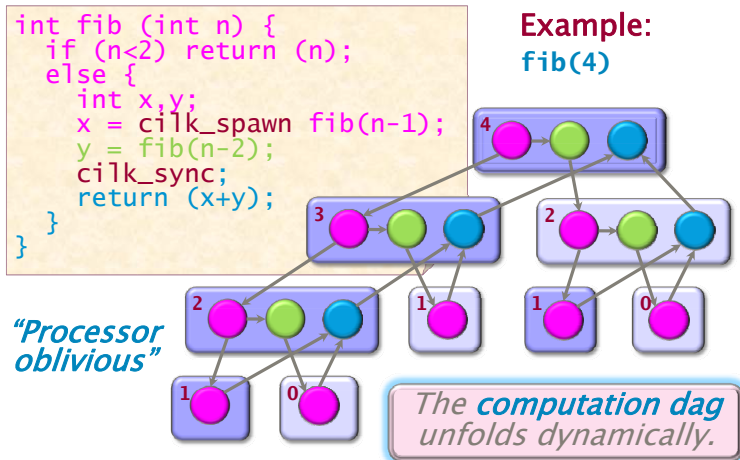
So does the (tuned) cache-oblivious matrix multiplication



Plan

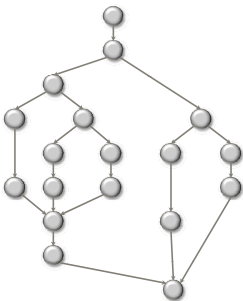
- 1 Data locality and cache misses
 - Hierarchical memories and their impact on our programs
 - Cache complexity and cache-oblivious algorithms put into practice
 - A detailed case study: counting sort
- 2 Multicore programming
 - Multicore architectures
 - Cilk / Cilk++ / Cilk Plus
 - The fork-join multithreaded programming model
- 3 GPU programming
 - The CUDA programming and memory models
 - Tiled matrix multiplication in CUDA
 - Optimizing Matrix Transpose with CUDA

The fork-join parallelism model



We shall also call this model **multithreaded parallelism**.

Work and span



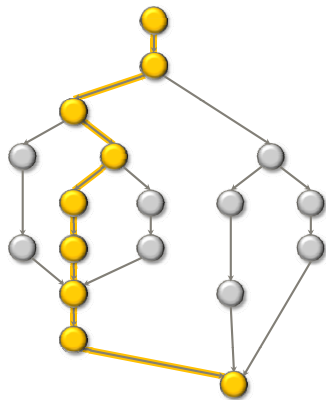
We define several performance measures. We assume an ideal situation: no cache issues, no interprocessor costs:

T_p is the minimum running time on p processors

T_1 is called the **work**, that is, the sum of the number of instructions at each node.

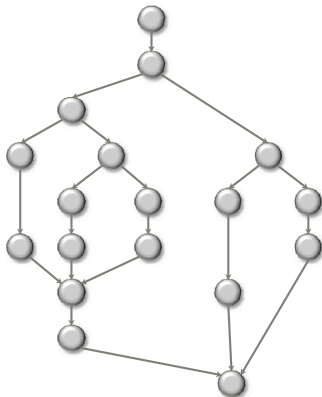
T_∞ is the minimum running time with infinitely many processors, called the **span**

The critical path length



Assuming all strands run in unit time, the longest path in the DAG is equal to T_∞ . For this reason, T_∞ is also referred to as the **critical path length**.

Work law

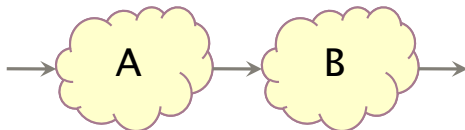


- We have: $T_p \geq T_1/p$.
- Indeed, in the best case, p processors can do p works per unit of time.

Speedup on p processors

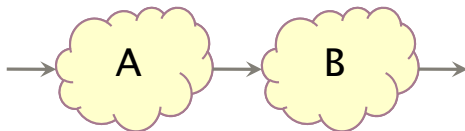
- T_1/T_p is called the **speedup on p processors**
- A parallel program execution can have:
 - **linear speedup**: $T_1/T_P = \Theta(p)$
 - **superlinear speedup**: $T_1/T_P = \omega(p)$ (not possible in this model, though it is possible in others)
 - **sublinear speedup**: $T_1/T_P = o(p)$

Series composition



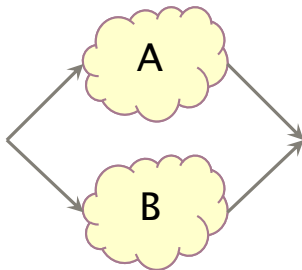
- Work?
- Span?

Series composition



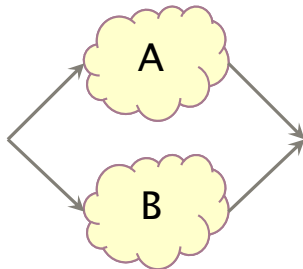
- Work: $T_1(A \cup B) = T_1(A) + T_1(B)$
- Span: $T_\infty(A \cup B) = T_\infty(A) + T_\infty(B)$

Parallel composition



- Work?
- Span?

Parallel composition



- Work: $T_1(A \cup B) = T_1(A) + T_1(B)$
- Span: $T_\infty(A \cup B) = \max(T_\infty(A), T_\infty(B))$

Some results in the fork-join parallelism model

Algorithm	Work	Span
Merge sort	$\Theta(n \lg n)$	$\Theta(\lg^3 n)$
Matrix multiplication	$\Theta(n^3)$	$\Theta(\lg n)$
Strassen	$\Theta(n^{\lg 7})$	$\Theta(\lg^2 n)$
LU-decomposition	$\Theta(n^3)$	$\Theta(n \lg n)$
Tableau construction	$\Theta(n^2)$	$\Omega(n^{\lg 3})$
FFT	$\Theta(n \lg n)$	$\Theta(\lg^2 n)$
Breadth-first search	$\Theta(E)$	$\Theta(d \lg V)$

We shall prove those results in the next lectures.

For loop parallelism in Cilk++

$$\begin{matrix} \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix} & \longrightarrow & \begin{pmatrix} a_{11} & a_{21} & \dots & a_{n1} \\ a_{12} & a_{22} & \dots & a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1n} & a_{2n} & \dots & a_{nn} \end{pmatrix} \\ A & & A^T \end{matrix}$$

```
cilk_for (int i=1; i<n; ++i) {
    for (int j=0; j<i; ++j) {
        double temp = A[i][j];
        A[i][j] = A[j][i];
        A[j][i] = temp;
    }
}
```

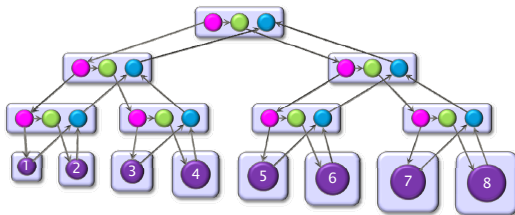
The iterations of a `cilk_for` loop execute in parallel.

Implementation of for loops in Cilk++

Up to details (next week!) the previous loop is compiled as follows, using a **divide-and-conquer implementation**:

```
void recur(int lo, int hi) {
    if (hi > lo) { // coarsen
        int mid = lo + (hi - lo)/2;
        cilk_spawn recur(lo, mid);
        recur(mid+1, hi);
        cilk_sync;
    } else
        for (int j=lo; j<hi; ++j) {
            double temp = A[hi][j];
            A[hi][j] = A[j][hi];
            A[j][hi] = temp;
        }
}
```

Analysis of parallel for loops



Here we do not assume that each strand runs in unit time.

- **Span of loop control:** $\Theta(\log(n))$
- **Max span of an iteration:** $\Theta(n)$
- **Span:** $\Theta(n)$
- **Work:** $\Theta(n^2)$
- **Parallelism:** $\Theta(n)$

For loops in the fork-join parallelism model: another example

```
cilk_for (int i = 1; i <= 8; i ++){
    f(i);
}
```

A *cilk_for* loop executes recursively as 2 for loops of $n/2$ iterations, adding a span of $\Theta(\log(n))$.

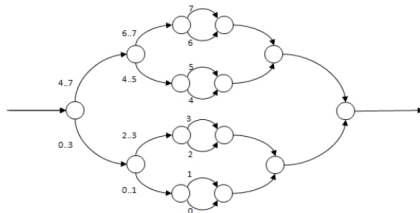
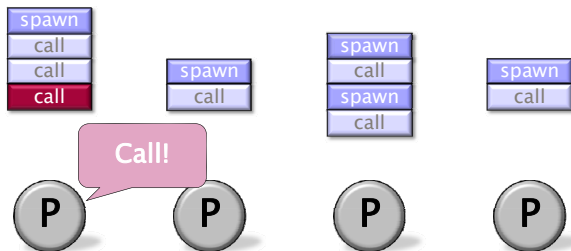
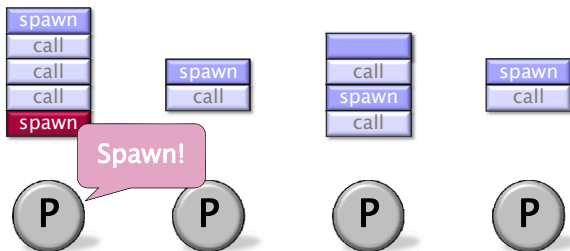


Figure: DAG for a *cilk_for* with 8 iterations.

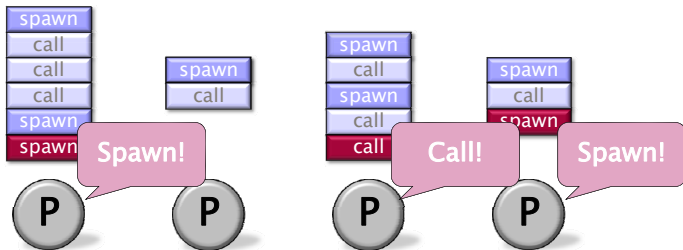
The work-stealing scheduler (1/11)



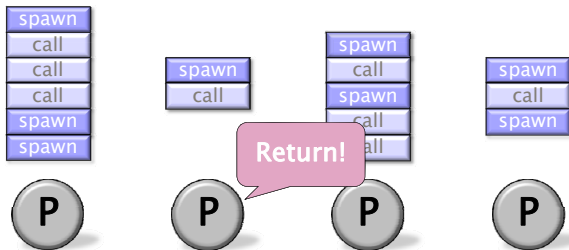
The work-stealing scheduler (2/11)



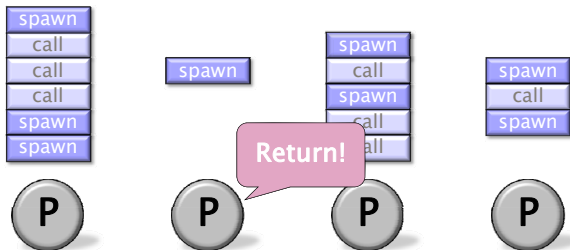
The work-stealing scheduler (3/11)



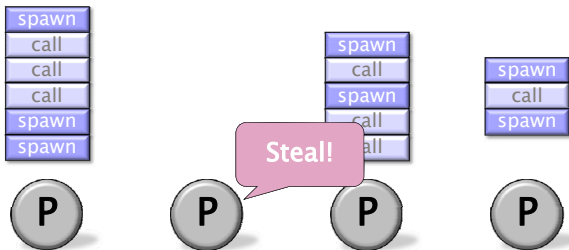
The work-stealing scheduler (4/11)



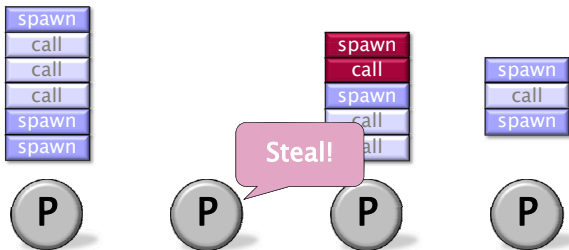
The work-stealing scheduler (5/11)



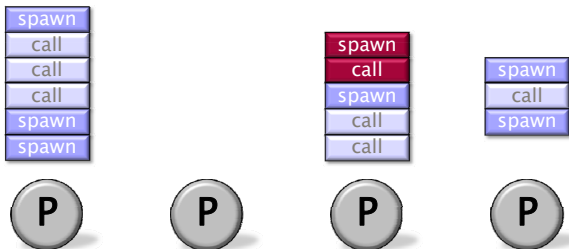
The work-stealing scheduler (6/11)



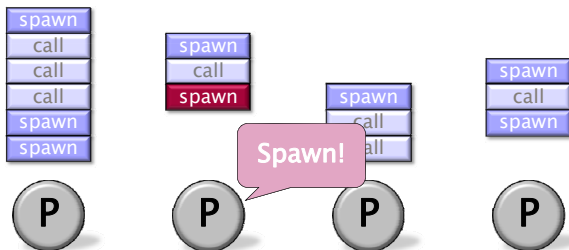
The work-stealing scheduler (7/11)



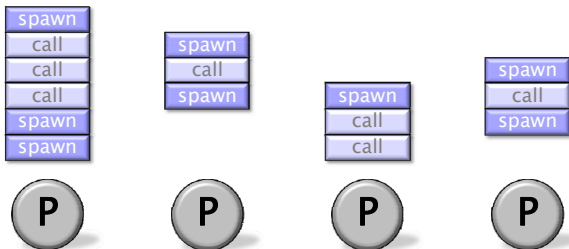
The work-stealing scheduler (8/11)



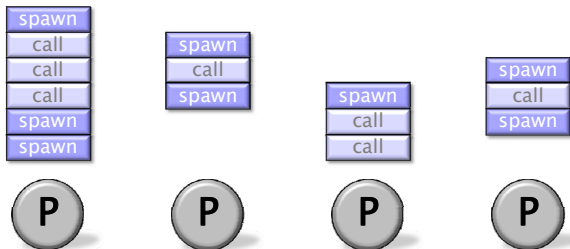
The work-stealing scheduler (9/11)



The work-stealing scheduler (10/11)



The work-stealing scheduler (11/11)



Performances of the work-stealing scheduler

Assume that

- each strand executes in unit time,
- for almost all “parallel steps” there are at least p strands to run,
- each processor is either working or stealing.

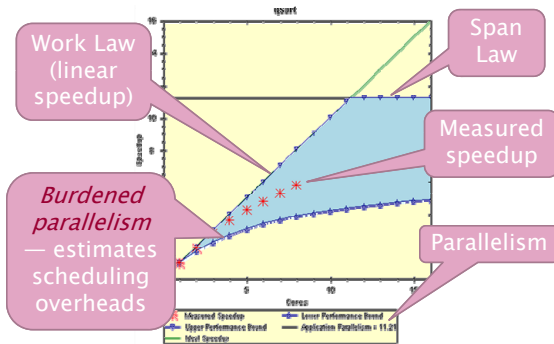
Then, the randomized work-stealing scheduler is expected to run in

$$T_P = T_1/p + O(T_\infty)$$

Overheads and burden

- Many factors (simplification assumptions of the fork-join parallelism model, architecture limitation, costs of executing the parallel constructs, overheads of scheduling) will make T_p larger in practice than $T_1/p + T_\infty$.
- One may want to estimate the impact of those factors:
 - ① by improving the estimate of the *randomized work-stealing complexity result*
 - ② by comparing a Cilk++ program with its C++ elision
 - ③ by estimating the costs of spawning and synchronizing
- Cilk++ estimates T_p as $T_p = T_1/p + 1.7 \text{ burden_span}$, where `burden_span` is 15000 instructions times the **number of continuation edges along the critical path**.

Cilkview



- **Cilkview** computes work and span to derive upper bounds on parallel performance
- **Cilkview** also estimates scheduling overhead to compute a burdened span for lower bounds.

The Fibonacci Cilk++ example

Code fragment

```
long fib(int n)
{
    if (n < 2) return n;
    long x, y;
    x = cilk_spawn fib(n-1);
    y = fib(n-2);
    cilk_sync;
    return x + y;
}
```

Fibonacci program timing

The environment for benchmarking:

- model name : Intel(R) Core(TM)2 Quad CPU Q6600 @ 2.40GHz
- L2 cache size : 4096 KB
- memory size : 3 GB

	#cores = 1	#cores = 2		#cores = 4	
n	timing(s)	timing(s)	speedup	timing(s)	speedup
30	0.086	0.046	1.870	0.025	3.440
35	0.776	0.436	1.780	0.206	3.767
40	8.931	4.842	1.844	2.399	3.723
45	105.263	54.017	1.949	27.200	3.870
50	1165.000	665.115	1.752	340.638	3.420

Quicksort

code in `cilk/examples/qsort`

```
void sample_qsor(int * begin, int * end)
{
    if (begin != end) {
        --end;
        int * middle = std::partition(begin, end,
            std::bind2nd(std::less<int>(), *end));
        using std::swap;
        swap(*end, *middle);
        cilk_spawn sample_qsor(begin, middle);
        sample_qsor(++middle, ++end);
        cilk_sync;
    }
}
```

Quicksort timing

Timing for sorting an array of integers:

	#cores = 1	#cores = 2		#cores = 4	
# of int	timing(s)	timing(s)	speedup	timing(s)	speedup
10×10^6	1.958	1.016	1.927	0.541	3.619
50×10^6	10.518	5.469	1.923	2.847	3.694
100×10^6	21.481	11.096	1.936	5.954	3.608
500×10^6	114.300	57.996	1.971	31.086	3.677

Matrix multiplication

[Code in cilk/examples/matrix](#)

Timing of multiplying a 687×837 matrix by a 837×1107 matrix

	iterative			recursive		
threshold	st(s)	pt(s)	su	st(s)	pt (s)	su
10	1.273	1.165	0.721	1.674	0.399	4.195
16	1.270	1.787	0.711	1.408	0.349	4.034
32	1.280	1.757	0.729	1.223	0.308	3.971
48	1.258	1.760	0.715	1.164	0.293	3.973
64	1.258	1.798	0.700	1.159	0.291	3.983
80	1.252	1.773	0.706	1.267	0.320	3.959

st = sequential time; pt = parallel time with 4 cores; su = speedup

The cilkview example from the documentation

Using `cilk_for` to perform operations over an array in parallel:

```
static const int COUNT = 4;
static const int ITERATION = 1000000;
long arr[COUNT];
long do_work(long k){
    long x = 15;
    static const int nn = 87;
    for (long i = 1; i < nn; ++i)
        x = x / i + k % i;
    return x;
}
int cilk_main(){
    for (int j = 0; j < ITERATION; j++)
        cilk_for (int i = 0; i < COUNT; i++)
            arr[i] += do_work( j * i + i + j);
}
```

1) Parallelism Profile

Work :	6,480,801,250 ins
Span :	2,116,801,250 ins
Burdened span :	31,920,801,250 ins
Parallelism :	3.06
Burdened parallelism :	0.20
Number of spawns/syncs:	3,000,000
Average instructions / strand :	720
Strands along span :	4,000,001
Average instructions / strand on span :	529

2) Speedup Estimate

2 processors:	0.21 - 2.00
4 processors:	0.15 - 3.06
8 processors:	0.13 - 3.06
16 processors:	0.13 - 3.06
32 processors:	0.12 - 3.06

A simple fix

Inverting the two for loops

```
int cilk_main()
{
    cilk_for (int i = 0; i < COUNT; i++)
        for (int j = 0; j < ITERATION; j++)
            arr[i] += do_work( j * i + i + j);
}
```

1) Parallelism Profile

Work :	5,295,801,529 ins
Span :	1,326,801,107 ins
Burdened span :	1,326,830,911 ins
Parallelism :	3.99
Burdened parallelism :	3.99
Number of spawns/syncs:	3
Average instructions / strand :	529,580,152
Strands along span :	5
Average instructions / strand on span:	265,360,221

2) Speedup Estimate

2 processors:	1.40 - 2.00
4 processors:	1.76 - 3.99
8 processors:	2.01 - 3.99
16 processors:	2.17 - 3.99
32 processors:	2.25 - 3.99

Timing

	#cores = 1	#cores = 2		#cores = 4	
version	timing(s)	timing(s)	speedup	timing(s)	speedup
original	7.719	9.611	0.803	10.758	0.718
improved	7.471	3.724	2.006	1.888	3.957

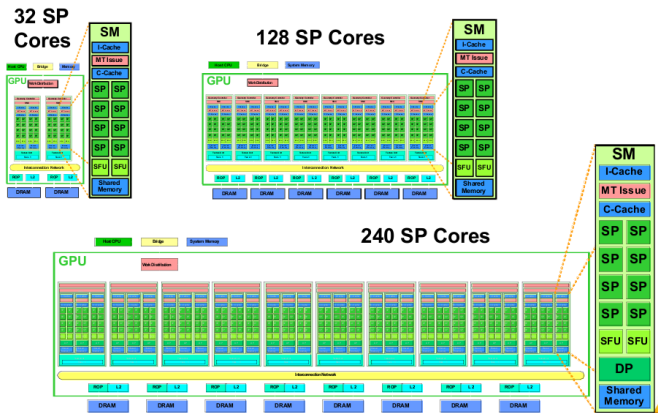
Summary and notes

Plan

- 1 Data locality and cache misses
 - Hierarchical memories and their impact on our programs
 - Cache complexity and cache-oblivious algorithms put into practice
 - A detailed case study: counting sort
- 2 Multicore programming
 - Multicore architectures
 - Cilk / Cilk++ / Cilk Plus
 - The fork-join multithreaded programming model
- 3 GPU programming
 - The CUDA programming and memory models
 - Tiled matrix multiplication in CUDA
 - Optimizing Matrix Transpose with CUDA

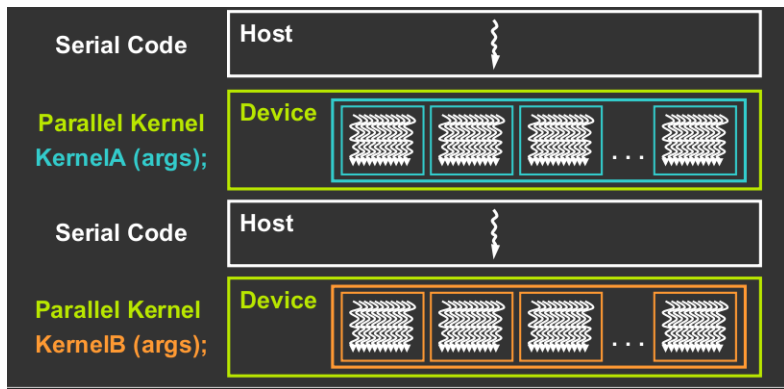
CUDA design goals

- Enable heterogeneous systems (i.e., CPU+GPU)
- Scale to 100's of cores, 1000's of parallel threads
- Use C/C++ with minimal extensions
- Let programmers focus on parallel algorithms



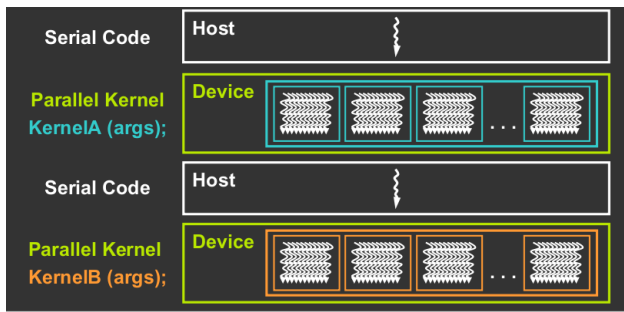
Heterogeneous programming (1/3)

- A CUDA program is a serial program with parallel kernels, all in C.
- The serial C code executes in a **host** (= CPU) thread
- The parallel kernel C code executes in many **device** threads across multiple GPU processing elements, called **streaming processors** (SP).



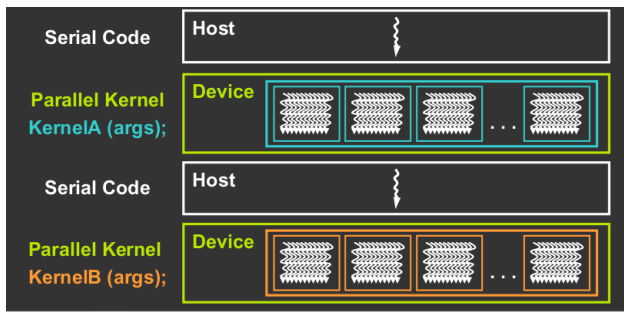
Heterogeneous programming (2/3)

- Thus, the parallel code (kernel) is launched and executed on a device by many threads.
- Threads are grouped into thread blocks.
- One kernel is executed at a time on the device.
- Many threads execute each kernel.



Heterogeneous programming (3/3)

- The parallel code is written for a thread
 - Each thread is free to execute a unique code path
 - Built-in **thread and block ID variables** are used to map each thread to a specific data tile (see next slide).
- Thus, each thread executes the same code on different data based on its thread and block ID.



Example: increment array elements (1/2)

Increment N-element vector a by scalar b



Let's assume $N=16$, $\text{blockDim}=4$ \rightarrow 4 blocks

```
int idx = blockDim.x * blockIdx.x + threadIdx.x;
```



$\text{blockIdx.x}=0$
 $\text{blockDim.x}=4$
 $\text{threadIdx.x}=0,1,2,3$
 $\text{idx}=0,1,2,3$



$\text{blockIdx.x}=1$
 $\text{blockDim.x}=4$
 $\text{threadIdx.x}=0,1,2,3$
 $\text{idx}=4,5,6,7$



$\text{blockIdx.x}=2$
 $\text{blockDim.x}=4$
 $\text{threadIdx.x}=0,1,2,3$
 $\text{idx}=8,9,10,11$



$\text{blockIdx.x}=3$
 $\text{blockDim.x}=4$
 $\text{threadIdx.x}=0,1,2,3$
 $\text{idx}=12,13,14,15$

See our example number 4 in `/usr/local/cs4402/examples/4`

Example: increment array elements (2/2)

CPU program

```
void increment_cpu(float *a, float b, int N)
{
    for (int idx = 0; idx < N; idx++)
        a[idx] = a[idx] + b;
}
```

```
void main()
{
    .....
    increment_cpu(a, b, N);
}
```

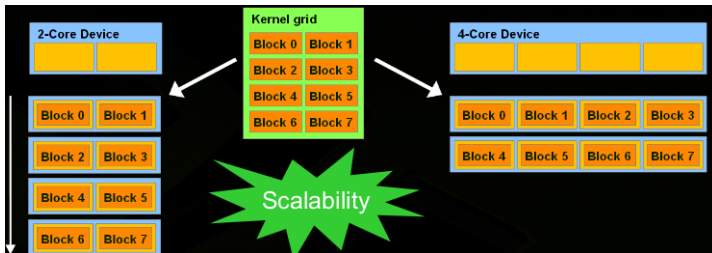
CUDA program

```
__global__ void increment_gpu(float *a, float b, int N)
{
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if (idx < N)
        a[idx] = a[idx] + b;
}
```

```
void main()
{
    .....
    dim3 dimBlock (blocksize);
    dim3 dimGrid( ceil( N / (float)blocksize) );
    increment_gpu<<<dimGrid, dimBlock>>>(a, b, N);
}
```

Thread blocks (1/2)

- A **Thread block** is a group of threads that can:
 - Synchronize their execution
 - Communicate via shared memory
- Within a grid, **thread blocks can run in any order**:
 - Concurrently or sequentially
 - Facilitates scaling of the same code across many devices



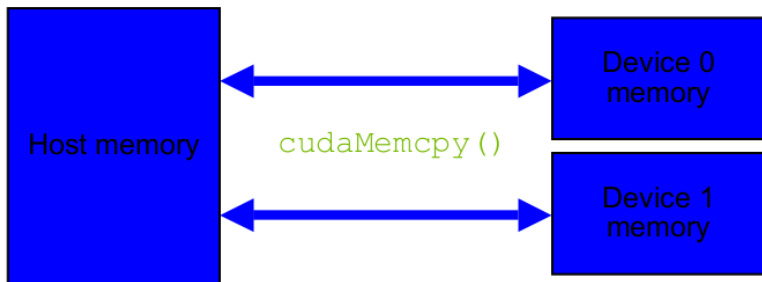
Thread blocks (2/2)

- Thus, within a grid, any possible interleaving of blocks must be valid.
- Thread blocks **may coordinate but not synchronize**
 - they may share pointers
 - they should not share locks (this can easily deadlock).
- The fact that thread blocks cannot synchronize gives **scalability**:
 - A kernel scales across any number of parallel cores
- However, within a thread block, threads may synchronize with barriers.
- That is, threads wait at the barrier until **all** threads in the **same block** reach the barrier.

Memory hierarchy (1/3)

Host (CPU) memory:

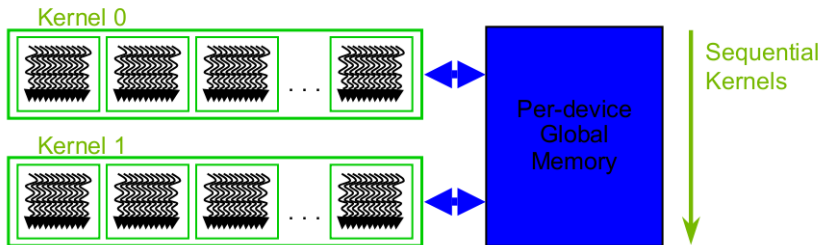
- Not directly accessible by CUDA threads



Memory hierarchy (2/3)

Global (on the device) memory:

- Also called **device memory**
- Accessible by all threads as well as host (CPU)
- Data lifetime = from allocation to deallocation



Memory hierarchy (3/3)

Shared memory:

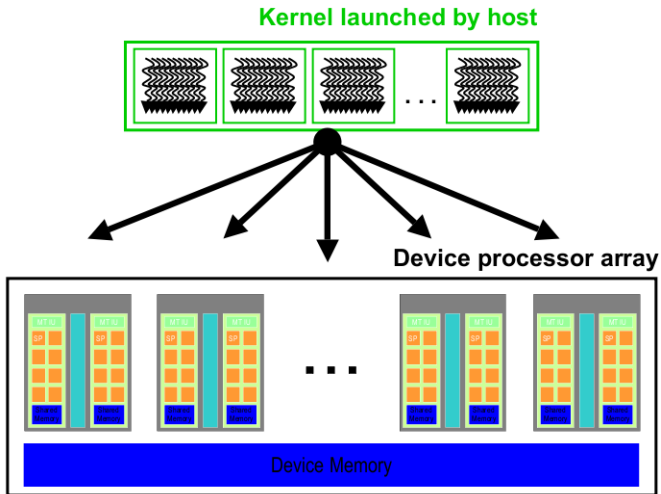
- Each thread block has its own shared memory, which is accessible only by the threads within that block
- Data lifetime = block lifetime

Local storage:

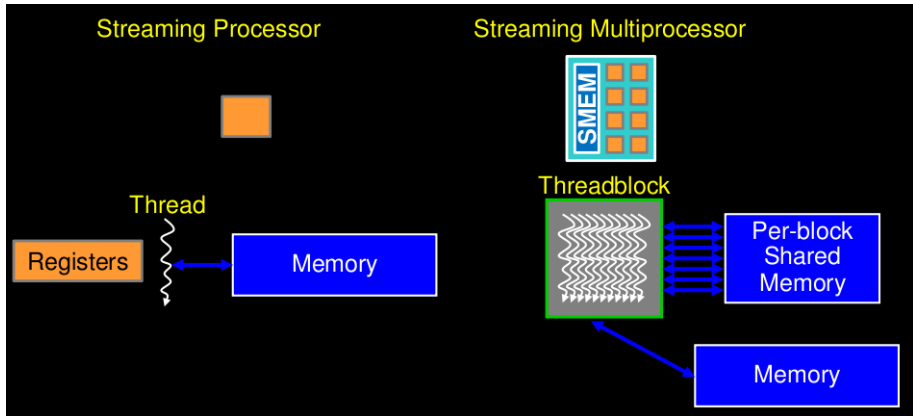
- Each thread has its own local storage
- Data lifetime = thread lifetime



Blocks run on multiprocessors



Streaming processors and multiprocessors



Hardware multithreading

- **Hardware allocates resources to blocks:**
 - blocks need: thread slots, registers, shared memory
 - blocks don't run until resources are available
- **Hardware schedules threads:**
 - threads have their own registers
 - any thread not waiting for something can run
 - context switching is free every cycle
- **Hardware relies on threads to hide latency:**
 - thus high parallelism is necessary for performance.



SIMT thread execution

- At each clock cycle, a multiprocessor executes the same instruction on a group of threads called a **warp**
 - The number of threads in a warp is the **warp size** (32 on G80)
 - A half-warp is the first or second half of a warp.
- Within a warp, threads
 - share instruction fetch/dispatch
 - some become inactive when code path diverges
 - hardware automatically handles divergence
- **Warps are the primitive unit of scheduling:**
 - each active block is split into warps in a well-defined way
 - threads within a warp are executed physically in parallel while warps and blocks are executed logically in parallel.



Plan

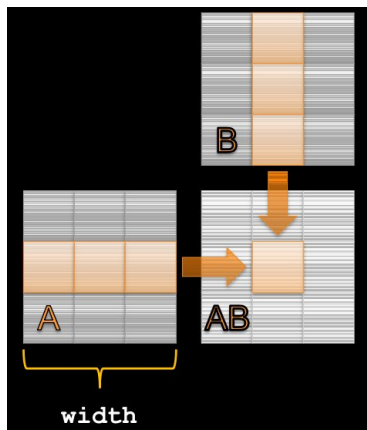
- 1 Data locality and cache misses
 - Hierarchical memories and their impact on our programs
 - Cache complexity and cache-oblivious algorithms put into practice
 - A detailed case study: counting sort
- 2 Multicore programming
 - Multicore architectures
 - Cilk / Cilk++ / Cilk Plus
 - The fork-join multithreaded programming model
- 3 GPU programming
 - The CUDA programming and memory models
 - Tiled matrix multiplication in CUDA
 - Optimizing Matrix Transpose with CUDA

Matrix multiplication (1/16)

- The goals of this example are:
 - Understanding how to write a kernel for a non-toy example
 - Understanding how to map work (and data) to the thread blocks
 - Understanding the importance of using shared memory
- We start by writing a naive kernel for matrix multiplication which does not use shared memory.
- Then we analyze the performance of this kernel and realize that it is limited by the global memory latency.
- Finally, we present a more efficient kernel, which takes advantage of a tile decomposition and makes use of shared memory.

Matrix multiplication (2/16)

- Consider multiplying two rectangular matrices A and B with respective formats $m \times n$ and $n \times p$. Define $C = A \times B$.
- Principle: each thread computes an element of C through a 2D grid with 2D thread blocks.



Matrix multiplication (3/16)

```
__global__ void mat_mul(float *a, float *b,
                        float *ab, int width)
{
    // calculate the row & col index of the element
    int row = blockIdx.y*blockDim.y + threadIdx.y;
    int col = blockIdx.x*blockDim.x + threadIdx.x;
    float result = 0;
    // do dot product between row of a and col of b
    for(int k = 0; k < width; ++k)
        result += a[row*width+k] * b[k*width+col];
    ab[row*width+col] = result;
}
```

Matrix multiplication (4/16)

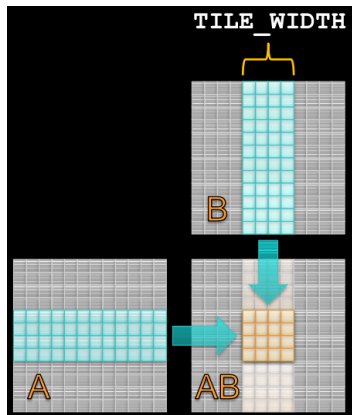
- Analyze the previous CUDA kernel for multiplying two rectangular matrices A and B with respective formats $m \times n$ and $n \times p$. Define $C = A \times B$.
- Each element of C is computed by one thread:
 - then each row of A is read p times and
 - each column of B is read m times, thus
 - **$2mnp$ reads in total for $2mnp$ flops.**
- Let t be an integer dividing m and p . We decompose C into $t \times t$ tiles. If tiles are computed one after another, then:
 - $(m/t)(tn)(p/t)$ slots are read in A
 - $(p/t)(tn)(m/t)$ slots are read in B , thus
 - **$2mnp/t$ reads in total for $2mnp$ flops.**
- For a CUDA implementation, $t = 16$ such that each tile is computed by one thread block.

Matrix multiplication (5/16)

- The previous explanation can be adapted to a particular GPU architecture, so as to estimate the performance of the first (naive) kernel.
- The first kernel has a **global memory access to flop ratio** (GMAC) of 8 Bytes / 2 ops, that is, 4 B/op.
- Suppose using a GeForce GTX 260, which has 805 GFLOPS peak performance.
- In order to reach **peak fp performance** we would need a memory bandwidth of $\text{GMAC} \times \text{Peak FLOPS} = 3.2 \text{ TB/s}$.
- Unfortunately, we only have 112 GB/s of actual **memory bandwidth** (BW) on a GeForce GTX 260.
- Therefore an upper bound on the performance of our implementation is $\text{BW} / \text{GMAC} = 28 \text{ GFLOPS}$.

Matrix multiplication (6/16)

- The picture below illustrates our second kernel
- Each thread block computes a tile in C , which is obtained as a dot product of tile-vector of A by a tile-vector of B .
- Tile size is chosen in order to maximize data locality.



Matrix multiplication (7/16)

- So a thread block computes a $t \times t$ tile of C .
- Each element in that tile is a dot-product of a row from A and a column from B .
- We view each of these dot-products as a sum of small dot products:

$$c_{i,j} = \sum_{k=0}^{t-1} a_{i,k} b_{k,j} + \sum_{k=t}^{2t-1} a_{i,k} b_{k,j} + \cdots + \sum_{k=n-1-t}^{n-1} a_{i,k} b_{k,j}$$

- Therefore we fix ℓ and then compute $\sum_{k=\ell t}^{(\ell+1)t-1} a_{i,k} b_{k,j}$ for all i, j in the working thread block.
- We do this for $\ell = 0, 1, \dots, (n/t - 1)$.
- This allows us to store the working tiles of A and B in shared memory.

Matrix multiplication (8/16)

- We assume that A , B , C are stored in row-major layout.
- Observe that for computing a tile in C our kernel code does need to know the number of rows in A .
- It just needs to know the **width** (number of columns) of A and B .

```
#define BLOCK_SIZE 16

    template <typename T>
__global__ void matrix_mul_ker(T* C, const T *A, const T *B,
    size_t wa, size_t wb)

// Block index; WARNING: should be at most 2^16 - 1
int bx = blockIdx.x; int by = blockIdx.y;

// Thread index
int tx = threadIdx.x; int ty = threadIdx.y;
```

Matrix multiplication (9/16)

- We need the position in `*A` of the first element of the first working tile from A ; we call it `aBegin`.
- We will need also the position in `*A` of the last element of the first working tile from A ; we call it `aEnd`.
- Moreover, we will need the offset between two consecutive working tiles of A ; we call it `aStep`.

```
int aBegin = wa * BLOCK_SIZE * by;
```

```
int aEnd = aBegin + wa - 1;
```

```
int aStep = BLOCK_SIZE;
```


Matrix multiplication (10/16)

- Similarly for B we have `bBegin` and `bStep`.
- We will not need a `bEnd` since once we are done with a row of A , we are also done with a column of B .
- Finally, we initialize the accumulator of the working thread; we call it `Csub`.

```
int bBegin = BLOCK_SIZE * bx;
```

```
int bStep = BLOCK_SIZE * wb;
```

```
int Csub = 0;
```

Matrix multiplication (11/16)

- The main loop starts by copying the working tiles of A and B to shared memory.

```
for(int a = aBegin, b = bBegin; a <= aEnd; a += aStep, b += bStep)
    // shared memory for the tile of A
    __shared__ int As[BLOCK_SIZE][BLOCK_SIZE];

    // shared memory for the tile of B
    __shared__ int Bs[BLOCK_SIZE][BLOCK_SIZE];

    // Load the tiles from global memory to shared memory
    // each thread loads one element of each tile
    As[ty][tx] = A[a + wa * ty + tx];
    Bs[ty][tx] = B[b + wb * ty + tx];

    // synchronize to make sure the matrices are loaded
    __syncthreads();
```

Matrix multiplication (12/16)

- Compute a small “dot-product” for each element in the working tile of C .

```
// Multiply the two tiles together
// each thread computes one element of the tile of C
for(int k = 0; k < BLOCK_SIZE; ++k) {
    Csub += As[ty][k] * Bs[k][tx];
}
// synchronize to make sure that the preceding computation
// done before loading two new tiles of A and B in the next iteration
__syncthreads();
}
```

Matrix multiplication (13/16)

- Once computed, the working tile of C is written to global memory.

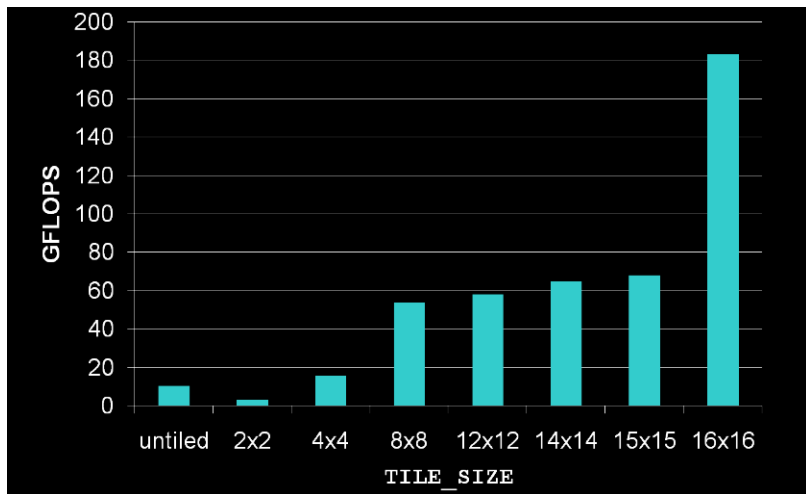
```
// Write the working tile of  $C$  to global memory;  
// each thread writes one element  
int c = wb * BLOCK_SIZE * by + BLOCK_SIZE * bx;  
C[c + wb * ty + tx] = Csub;
```

Matrix multiplication (14/16)

- Each thread block should have many threads:
 - `TILE_WIDTH = 16` implies $16 \times 16 = 256$ threads
- There should be many thread blocks:
 - A 1024×1024 matrix would require 4096 thread blocks.
 - Since one streaming multiprocessor (SM) can handle 768 threads, each SM will process 3 thread blocks, leading it **full occupancy**.
- Each thread block performs 2×256 reads of a 4-byte float while performing $256 \times (2 \times 16) = 8,192$ fp ops:
 - Memory bandwidth is no longer limiting factor

Matrix multiplication (15/16)

- Experimentation performed on a GT200.
- **Tiling** and using **shared memory** were clearly worth the effort.



Matrix multiplication (16/16)

- Effective use of different memory resources reduces the number of accesses to global memory
- But these resources are finite!
- The more memory locations each thread requires, the fewer threads an SM can accommodate.

Resource	Per GT200 SM	Full Occupancy on GT200
Registers	16384	$\leq 16384 / 768$ threads = 21 per thread
<u>shared</u> Memory	16KB	$\leq 16\text{KB} / 8$ blocks = 2KB per block

Plan

- 1 Data locality and cache misses
 - Hierarchical memories and their impact on our programs
 - Cache complexity and cache-oblivious algorithms put into practice
 - A detailed case study: counting sort
- 2 Multicore programming
 - Multicore architectures
 - Cilk / Cilk++ / Cilk Plus
 - The fork-join multithreaded programming model
- 3 GPU programming
 - The CUDA programming and memory models
 - Tiled matrix multiplication in CUDA
 - Optimizing Matrix Transpose with CUDA

Matrix transpose characteristics (1/2)

- We optimize a transposition code for a matrix of floats. This operates out-of-place:
 - input and output matrices address separate memory locations.
- For simplicity, we consider an $n \times n$ matrix where 32 divides n .
- We focus on the device code:
 - the host code performs typical tasks: data allocation and transfer between host and device, the launching and timing of several kernels, result validation, and the deallocation of host and device memory.
- Benchmarks illustrate this section:
 - we compare our **matrix transpose** kernels against a **matrix copy** kernel,
 - for each kernel, we compute the **effective bandwidth**, calculated in GB/s as twice the size of the matrix (once for reading the matrix and once for writing) divided by the time of execution,
 - Each operation is run NUM_REFS times (for **normalizing the measurements**),
 - This looping is performed **once over the kernel** and once **within the kernel**,
 - The difference between these two timings is kernel launch and synchronization overheads.

Matrix transpose characteristics (2/2)

- We present hereafter different kernels called from the host code, each addressing different performance issues.
- All kernels in this study launch thread blocks of dimension 32×8 , where each block transposes (or copies) a tile of dimension 32×32 .
- As such, the parameters `TILE_DIM` and `BLOCK_ROWS` are set to 32 and 8, respectively.
- Using a thread block with fewer threads than elements in a tile is advantageous for the matrix transpose:
 - each thread transposes several matrix elements, four in our case, and much of the cost of calculating the indices is amortized over these elements.
- This study is based on a technical report by Greg Ruetsch (NVIDIA) and Paulius Micikevicius (NVIDIA).

A simple copy kernel (1/2)

```
__global__ void copy(float *odata, float* idata, int width,
                    int height, int nreps)
{
    int xIndex = blockIdx.x*TILE_DIM + threadIdx.x;
    int yIndex = blockIdx.y*TILE_DIM + threadIdx.y;
    int index  = xIndex + width*yIndex;

    for (int r=0; r < nreps; r++) { // normalization outer loop
        for (int i=0; i<TILE_DIM; i+=BLOCK_ROWS) {
            odata[index+i*width] = idata[index+i*width];
        }
    }
}
```

A simple copy kernel (2/2)

- `odata` and `idata` are pointers to the input and output matrices,
- `width` and `height` are the matrix x and y dimensions,
- `nreps` determines how many times the loop over data movement between matrices is performed.
- In this kernel, `xIndex` and `yIndex` are global 2D matrix indices,
- used to calculate `index`, the 1D index used to access matrix elements.

```
__global__ void copy(float *odata, float* idata, int width,
                    int height, int nreps)
{
    int xIndex = blockIdx.x*TILE_DIM + threadIdx.x;
    int yIndex = blockIdx.y*TILE_DIM + threadIdx.y;
    int index  = xIndex + width*yIndex;

    for (int r=0; r < nreps; r++) {
        for (int i=0; i<TILE_DIM; i+=BLOCK_ROWS) {
            odata[index+i*width] = idata[index+i*width];
        } } }
```

A naive transpose kernel

```
_global__ void transposeNaive(float *odata, float* idata,
                             int width, int height, int nreps)
{
    int xIndex = blockIdx.x*TILE_DIM + threadIdx.x;
    int yIndex = blockIdx.y*TILE_DIM + threadIdx.y;
    int index_in = xIndex + width * yIndex;
    int index_out = yIndex + height * xIndex;
    for (int r=0; r < nreps; r++) {
        for (int i=0; i<TILE_DIM; i+=BLOCK_ROWS) {
            odata[index_out+i] = idata[index_in+i*width];
        }
    }
}
```

Naive transpose kernel vs copy kernel

The performance of these two kernels on a 2048x2048 matrix using a GTX280 is given in the following table:

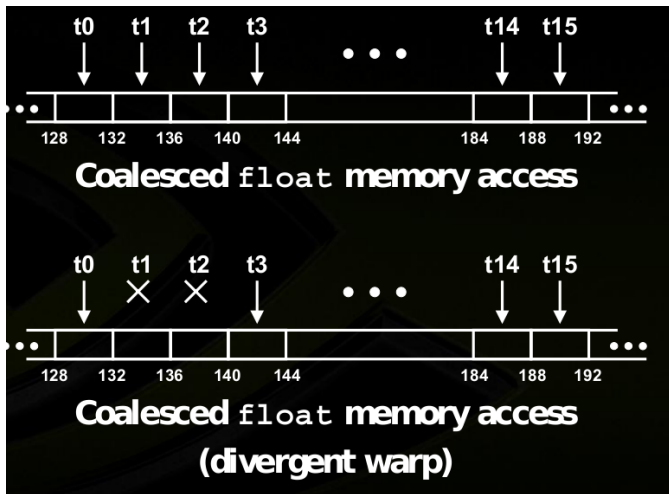
	Effective Bandwidth (GB/s) 2048x2048, GTX 280	
	Loop over kernel	Loop in kernel
Simple Copy	96.9	81.6
Naïve Transpose	2.2	2.2

The minor differences in code between the copy and naive transpose kernels have a profound effect on performance.

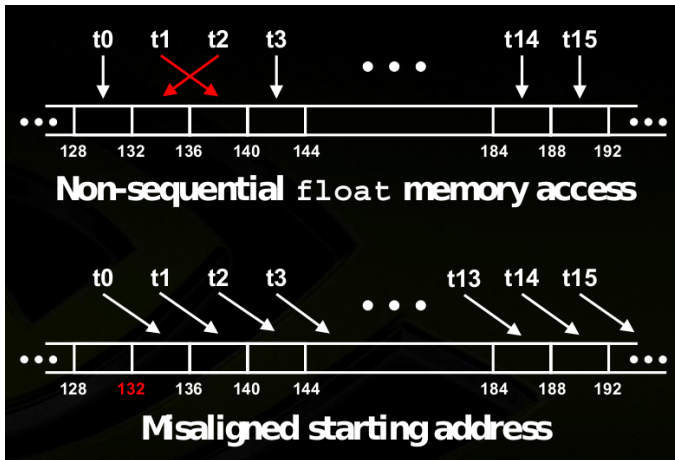
Coalesced Transpose (1/11)

- Because device memory has a much higher latency and lower bandwidth than on-chip memory, special attention must be paid to: **how global memory accesses are performed?**
- The simultaneous global memory accesses by each thread of a half-warp (16 threads on G80) during the execution of a single read or write instruction will be **coalesced** into a single access if:
 - ① The size of the memory element accessed by each thread is either 4, 8, or 16 bytes.
 - ② The address of the first element is aligned to 16 times the element's size.
 - ③ The elements form a contiguous block of memory.
 - ④ The i -th element is accessed by the i -th thread in the half-warp.
- Last two requirements are relaxed with compute capabilities of 1.2.
- Coalescing happens even if some threads do not access memory (**divergent warp**)

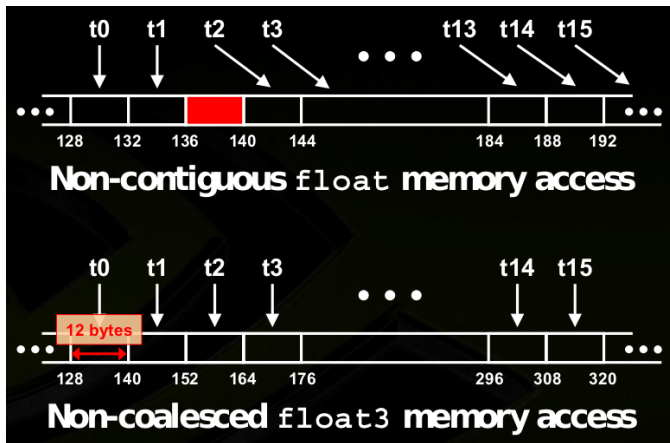
Coalesced Transpose (2/11)



Coalesced Transpose (3/11)



Coalesced Transpose (4/11)



Coalesced Transpose (5/11)

- **Allocating device memory through `cudaMalloc()`** and choosing **`TILE_DIM` to be a multiple of 16 ensures alignment** with a segment of memory, therefore all loads from `idata` are coalesced.
- Coalescing behavior differs between the simple copy and naive transpose kernels when writing to `odata`.
- In the case of the naive transpose, for each iteration of the `i`-loop a half warp writes one half of a column of floats to different segments of memory:
 - resulting in 16 separate memory transactions,
 - regardless of the compute capability.

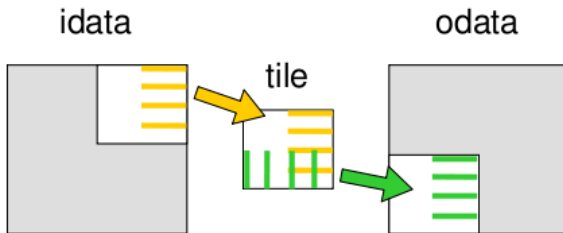
Coalesced Transpose (6/11)

- The way to avoid uncoalesced global memory access is
 - ① to read the data into shared memory and,
 - ② have each half warp access non-contiguous locations in shared memory in order to write contiguous data to odata.
- There is no performance penalty for non-contiguous access patterns in shared memory as there is in global memory.
- a `__syncthreads()` call is required to ensure that all reads from `idata` to shared memory have completed before writes from shared memory to `odata` commence.

Coalesced Transpose (7/11)

```
__global__ void transposeCoalesced(float *odata,
                                   float *idata, int width, int height) // no nreps param
{
    __shared__ float tile[TILE_DIM][TILE_DIM];
    int xIndex = blockIdx.x*TILE_DIM + threadIdx.x;
    int yIndex = blockIdx.y*TILE_DIM + threadIdx.y;
    int index_in = xIndex + (yIndex)*width;
    xIndex = blockIdx.y * TILE_DIM + threadIdx.x;
    yIndex = blockIdx.x * TILE_DIM + threadIdx.y;
    int index_out = xIndex + (yIndex)*height;
    for (int i=0; i<TILE_DIM; i+=BLOCK_ROWS) {
        tile[threadIdx.y+i][threadIdx.x] =
            idata[index_in+i*width];
    }
    __syncthreads();
    for (int i=0; i<TILE_DIM; i+=BLOCK_ROWS) {
        odata[index_out+i*height] =
            tile[threadIdx.x][threadIdx.y+i];
    }
}
```

Coalesced Transpose (8/11)



- 1 The half warp writes four half rows of the `idata` matrix tile to the shared memory `32x32` array `tile` indicated by the yellow line segments.
- 2 After a `__syncthreads()` call to ensure all writes to `tile` are completed,
- 3 the half warp writes four half columns of `tile` to four half rows of an `odata` matrix tile, indicated by the green line segments.

Coalesced Transpose (9/11)

	Effective Bandwidth (GB/s) 2048x2048, GTX 280	
	Loop over kernel	Loop in kernel
Simple Copy	96.9	81.6
Naïve Transpose	2.2	2.2
Coalesced Transpose	16.5	17.1

While there is a dramatic increase in effective bandwidth of the coalesced transpose over the naïve transpose, there still remains a large performance gap between the coalesced transpose and the copy:

- One possible cause of this performance gap could be the synchronization barrier required in the coalesced transpose.
- This can be easily assessed using the following copy kernel which utilizes shared memory and contains a `__syncthreads()` call.

Coalesced Transpose (10/11)

```
__global__ void copySharedMem(float *odata, float *idata,
                              int width, int height) // no nreps param
{
    __shared__ float tile[TILE_DIM][TILE_DIM];
    int xIndex = blockIdx.x*TILE_DIM + threadIdx.x;
    int yIndex = blockIdx.y*TILE_DIM + threadIdx.y;
    int index = xIndex + width*yIndex;
    for (int i=0; i<TILE_DIM; i+=BLOCK_ROWS) {
        tile[threadIdx.y+i][threadIdx.x] =
            idata[index+i*width];
    }
    __syncthreads();
    for (int i=0; i<TILE_DIM; i+=BLOCK_ROWS) {
        odata[index+i*width] =
            tile[threadIdx.y+i][threadIdx.x];
    }
}
```


Coalesced Transpose (11/11)

	Effective Bandwidth (GB/s) 2048x2048, GTX 280	
	Loop over kernel	Loop in kernel
Simple Copy	96.9	81.6
Shared Memory Copy	80.9	81.1
Naïve Transpose	2.2	2.2
Coalesced Transpose	16.5	17.1

The shared memory copy results seem to suggest that the use of shared memory with a synchronization barrier has little effect on the performance, certainly as far as the *Loop in kernel* column indicates when comparing the simple copy and shared memory copy.

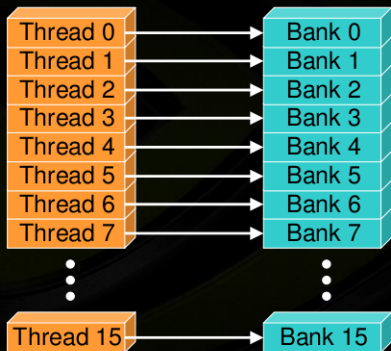
Shared memory bank conflicts (1/6)

- 1 Shared memory is divided into 16 equally-sized memory modules, called **banks**, which are organized such that successive 32-bit words are assigned to successive banks.
- 2 These banks can be accessed simultaneously, and to achieve maximum bandwidth to and from shared memory the **threads in a half warp should access shared memory associated with different banks.**
- 3 The **exception to this rule is** when all threads in a half warp read the same shared memory address, which results in a broadcast where the data at that address is sent to all threads of the half warp in one transaction.
- 4 One can use the `warp_serialize` flag when profiling CUDA applications to determine whether shared memory bank conflicts occur in any kernel.

Shared memory bank conflicts (2/6)

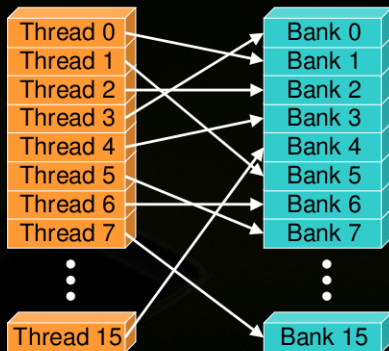
- No bank conflicts

- Linear addressing
stride == 1



- No bank conflicts

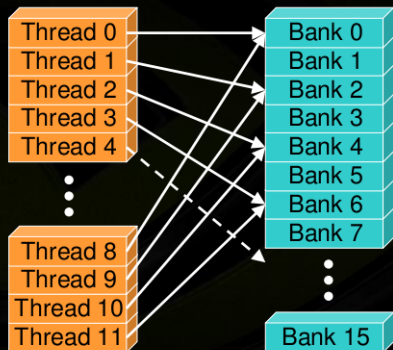
- Random 1:1 permutation



Shared memory bank conflicts (3/6)

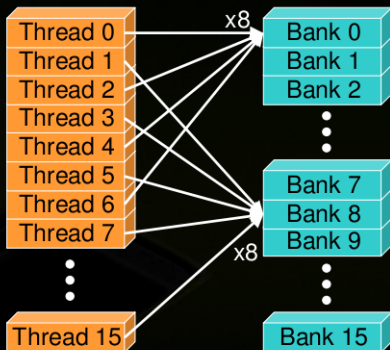
2-way bank conflicts

- Linear addressing stride == 2



8-way bank conflicts

- Linear addressing stride == 8



Shared memory bank conflicts (4/6)

- 1 The coalesced transpose uses a 32×32 shared memory array of floats.
- 2 For this sized array, all data in columns k and $k+16$ are mapped to the same bank.
- 3 As a result, when writing partial columns from `tile` in shared memory to rows in `odata` the half warp experiences a 16-way bank conflict and serializes the request.
- 4 A simple way to avoid this conflict is to pad the shared memory array by one column:

```
__shared__ float tile[TILE_DIM][TILE_DIM+1];
```

Shared memory bank conflicts (5/6)

- The padding does not affect shared memory bank access pattern when writing a half warp to shared memory, which remains conflict free,
- but by adding a single column now the access of a half warp of data in a column is also conflict free.
- The performance of the kernel, now coalesced and memory bank conflict free, is added to our table on the next slide.

Shared memory bank conflicts (6/6)

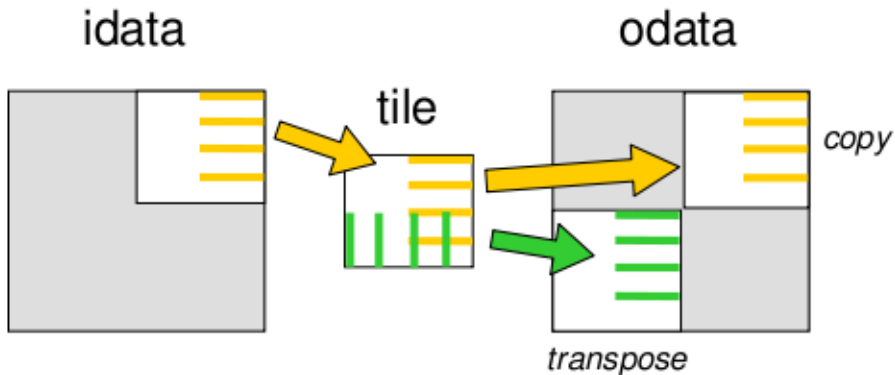
	Effective Bandwidth (GB/s) 2048x2048, GTX 280	
	Loop over kernel	Loop in kernel
Simple Copy	96.9	81.6
Shared Memory Copy	80.9	81.1
Naïve Transpose	2.2	2.2
Coalesced Transpose	16.5	17.1
Bank Conflict Free Transpose	16.6	17.2

- While padding the shared memory array did eliminate shared memory bank conflicts, as was confirmed by checking the `warp_serialize` flag with the CUDA profiler, it has little effect (when implemented at this stage) on performance.
- As a result, there is still a large performance gap between the coalesced and shared memory bank conflict free transpose and the shared memory copy.

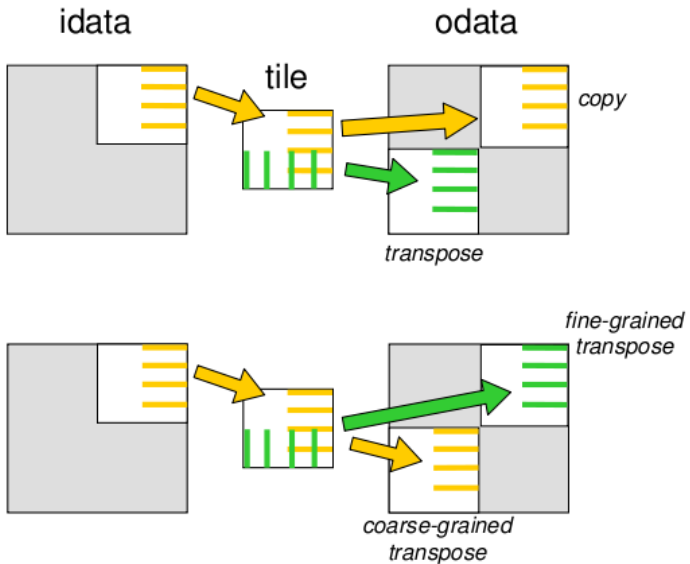
Decomposing Transpose (1/6)

- To investigate further, we revisit the data flow for the transpose and compare it to that of the copy.
- There are essentially two differences between the copy code and the transpose:
 - transposing the data within a tile, and
 - writing data to transposed tile.
- We can isolate the performance between each of these two components by implementing two kernels that individually perform just one of these components:
 - fine-grained transpose:** this kernel transposes the data within a tile, but writes the tile to the location.
 - coarse-grained transpose:** this kernel writes the tile to the transposed location in the odata matrix, but does not transpose the data within the tile.

Decomposing Transpose (2/6)



Decomposing Transpose (3/6)



Decomposing Transpose (4/6)

```
_global__ void transposeFineGrained(float *odata,
    float *idata, int width, int height)
{
    __shared__ float block[TILE_DIM][TILE_DIM+1];
    int xIndex = blockIdx.x * TILE_DIM + threadIdx.x;
    int yIndex = blockIdx.y * TILE_DIM + threadIdx.y;
    int index = xIndex + (yIndex)*width;
    for (int i=0; i < TILE_DIM; i += BLOCK_ROWS) {
        block[threadIdx.y+i][threadIdx.x] =
            idata[index+i*width];
    }
    __syncthreads();
    for (int i=0; i < TILE_DIM; i += BLOCK_ROWS) {
        odata[index+i*height] =
            block[threadIdx.x][threadIdx.y+i];
    }
}
```

Decomposing Transpose (5/6)

```
__global__ void transposeCoarseGrained(float *odata,
    float *idata, int width, int height)
{
    __shared__ float block[TILE_DIM][TILE_DIM+1];
    int xIndex = blockIdx.x * TILE_DIM + threadIdx.x;
    int yIndex = blockIdx.y * TILE_DIM + threadIdx.y;
    int index_in = xIndex + (yIndex)*width;
    xIndex = blockIdx.y * TILE_DIM + threadIdx.x;
    yIndex = blockIdx.x * TILE_DIM + threadIdx.y;
    int index_out = xIndex + (yIndex)*height;
    for (int i=0; i<TILE_DIM; i += BLOCK_ROWS) {
        block[threadIdx.y+i][threadIdx.x] =
            idata[index_in+i*width];
    } syncthreads();
    for (int i=0; i<TILE_DIM; i += BLOCK_ROWS) {
        odata[index_out+i*height] =
            block[threadIdx.y+i][threadIdx.x];
    } }
```

Decomposing Transpose (6/6)

	Effective Bandwidth (GB/s) 2048x2048, GTX 280	
	Loop over kernel	Loop in kernel
Simple Copy	96.9	81.6
Shared Memory Copy	80.9	81.1
Naïve Transpose	2.2	2.2
Coalesced Transpose	16.5	17.1
Bank Conflict Free Transpose	16.6	17.2
<i>Fine-grained Transpose</i>	<i>80.4</i>	<i>81.5</i>
<i>Coarse-grained Transpose</i>	<i>16.7</i>	<i>17.1</i>

The fine-grained transpose has performance similar to the shared memory copy, whereas the coarse-grained transpose has roughly the performance of the coalesced transpose. Thus the performance bottleneck lies in writing data to the transposed location in global memory.

Partition Camping (1/4)

- Just as shared memory performance can be degraded via bank conflicts, an analogous performance degradation can occur with global memory access through **partition camping**.
- Global memory is divided into either 6 **partitions** (on 8- and 9-series GPUs) or 8 partitions (on 200- and 10-series GPUs) of 256-byte width.
- To use global memory effectively, concurrent accesses to global memory by all active warps should be divided evenly amongst partitions.
- **partition camping** occurs when:
 - global memory accesses are directed through a subset of partitions,
 - causing requests to queue up at some partitions while other partitions go unused.

Partition Camping (2/4)

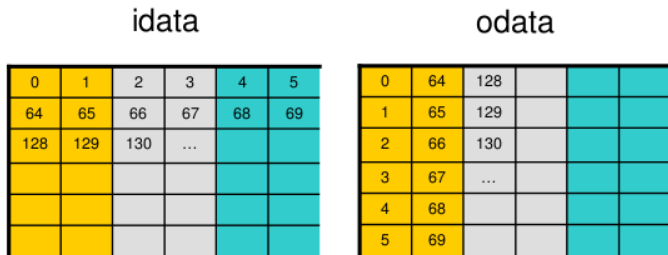
- Since partition camping concerns how active thread blocks behave, the issue of how thread blocks are scheduled on multiprocessors is important.
- When a kernel is launched, the order in which blocks are assigned to multiprocessors is determined by the one-dimensional block ID defined as:

```
bid = blockIdx.x + blockDim.x*blockIdx.y;
```

which is a row-major ordering of the blocks in the grid.

- Once maximum occupancy is reached, additional blocks are assigned to multiprocessors as needed.
- How quickly and the order in which blocks complete cannot be determined.
- So active blocks are initially contiguous but become less contiguous as execution of the kernel progresses.

Partition Camping (3/4)



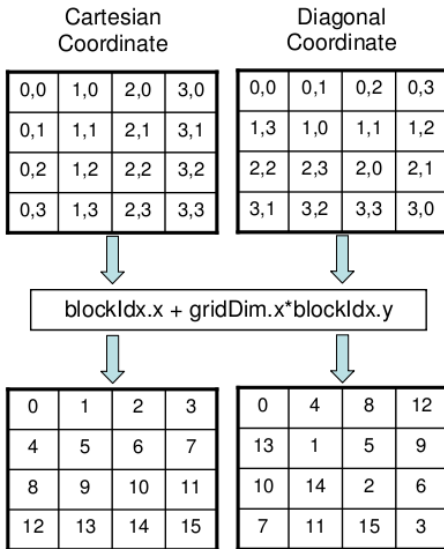
- With 8 partitions of 256-byte width, all data in strides of 2048 bytes (or 512 floats) map to the same partition.
- Any float matrix with $512 \times k$ columns, such as our 2048x2048 matrix, will contain columns whose elements map to a single partition.
- With tiles of 32×32 floats whose one-dimensional block IDs are shown in the figures, the mapping of `idata` and `odata` onto the partitions is depicted below.

Partition Camping (4/4)

idata						odata					
0	1	2	3	4	5	0	64	128			
64	65	66	67	68	69	1	65	129			
128	129	130	...			2	66	130			
						3	67	...			
						4	68				
						5	69				

- Concurrent blocks will be accessing tiles row-wise in `idata` which will be roughly equally distributed amongst partitions
- However these blocks will access tiles column-wise in `odata` which will typically access global memory through just a few partitions.
- Just as with shared memory, padding would be an option (potentially expensive) but there is a better one ...

Diagonal block reordering (1/7)



Diagonal block reordering (2/7)

- The key idea is to view the grid under a **diagonal coordinate system**.
- If `blockIdx.x` and `blockIdx.y` represent the diagonal coordinates, then (for block-square matrices) the corresponding Cartesian coordinates are given by the following mapping:

```
blockIdx_y = blockIdx.x;
```

```
blockIdx_x = (blockIdx.x+blockIdx.y)%gridDim.x;
```

- One would simply include the previous two lines of code at the beginning of the kernel, and write the kernel assuming the Cartesian interpretation of `blockIdx` fields, except using `blockIdx_x` and `blockIdx_y` in place of `blockIdx.x` and `blockIdx.y`, respectively, throughout the kernel.
- This is precisely what is done in the `transposeDiagonal` kernel hereafter.

Decomposing Transpose (3/7)

```
__global__ void transposeDiagonal(float *odata,
    float *idata, int width, int height)
{
    __shared__ float tile[TILE_DIM][TILE_DIM+1];
    int blockIdx_x, blockIdx_y;
    // diagonal reordering
    if (width == height) {
        blockIdx_y = blockIdx.x;
        blockIdx_x = (blockIdx.x+blockIdx.y)%gridDim.x;
    } else {
        int bid = blockIdx.x + gridDim.x*blockIdx.y;
        blockIdx_y = bid%gridDim.y;
        blockIdx_x = ((bid/gridDim.y)+blockIdx_y)%gridDim.x;
    }
}
```

Decomposing Transpose (4/7)

```
int xIndex = blockIdx_x*TILE_DIM + threadIdx.x;
int yIndex = blockIdx_y*TILE_DIM + threadIdx.y;
int index_in = xIndex + (yIndex)*width;
xIndex = blockIdx_y*TILE_DIM + threadIdx.x;
yIndex = blockIdx_x*TILE_DIM + threadIdx.y;
int index_out = xIndex + (yIndex)*height;
    for (int i=0; i<TILE_DIM; i+=BLOCK_ROWS) {
        tile[threadIdx.y+i][threadIdx.x] =
            idata[index_in+i*width];
    }
    __syncthreads();
    for (int i=0; i<TILE_DIM; i+=BLOCK_ROWS) {
        odata[index_out+i*height] =
            tile[threadIdx.x][threadIdx.y+i];
    }
}
```

Diagonal block reordering (5/7)

idata

0	1	2	3	4	5
64	65	66	67	68	69
128	129	130	...		

Cartesian

odata

0	64	128			
1	65	129			
2	66	130			
3	67	...			
4	68				
5	69				

Diagonal

0	64	128			
	1	65	129		
		2	66	130	
			3	67	...
				4	68
					5

0					
64	1				
128	65	2			
	129	66	3		
		130	67	4	
			...	68	5

Diagonal block reordering (6/7)

	Effective Bandwidth (GB/s) 2048x2048, GTX 280	
	Loop over kernel	Loop in kernel
Simple Copy	96.9	81.6
Shared Memory Copy	80.9	81.1
Naïve Transpose	2.2	2.2
Coalesced Transpose	16.5	17.1
Bank Conflict Free Transpose	16.6	17.2
<i>Fine-grained Transpose</i>	80.4	81.5
<i>Coarse-grained Transpose</i>	16.7	17.1
Diagonal	69.5	78.3

Diagonal block reordering (7/7)

- The bandwidth measured when looping within the kernel over the read and writes to global memory is within a few percent of the shared memory copy.
- When looping over the kernel, the performance degrades slightly, likely due to additional computation involved in calculating `blockIdx_x` and `blockIdx_y`. However, even with this performance degradation the diagonal transpose has over four times the bandwidth of the other complete transposes.

	Effective Bandwidth (GB/s) 2048x2048, GTX 280	
	Loop over kernel	Loop in kernel
Simple Copy	96.9	81.6
Shared Memory Copy	80.9	81.1
Naïve Transpose	2.2	2.2
Coalesced Transpose	16.5	17.1
Bank Conflict Free Transpose	16.6	17.2
<i>Fine-grained Transpose</i>	80.4	81.5
<i>Coarse-grained Transpose</i>	16.7	17.1
Diagonal	69.5	78.3

Summary and notes