Plan



	< 🗆	▶ ★@ ▶ ★ 臣 ▶ ★ 臣 ▶ 二 臣	$\mathcal{O}\mathcal{A}\mathcal{O}$		< □ >	★@▶★ 문▶ ★ 문▶ _ 문	うくで
(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	1 / 60	(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	2 / 60
	Multi-core Architecture				Multi-core Architecture Multi-core processor		

Plan

Multi-core Architecture

- Multi-core processor
- CPU Cache
- CPU Coherence

Concurrency Platforms

- PThreads
- TBB
- Open MP
- Cilk ++
- Race Conditions and Cilkscreen
- MMM in Cilk++



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで Introduction to Multicore Programming CS 4435 - CS 9624

(Moreno Maza)

3 / 60

Introduction to Multicore Programming



Multi-core Architecture Multi-core processor



	< 🗆	▷ ▲@ ▶ ▲글 ▶ ▲글 ▶ _ 글	$\mathcal{O}\mathcal{A}\mathcal{O}$		< □)	・ 本間 と 本語 と 不語 と 三語	うくつ
(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	5 / 60	(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	6 / 60
N	Aulti-core Architecture Multi-core processor			N	Iulti-core Architecture Multi-core processor		

Multi-core processor

- A multi-core processor is an integrated circuit to which two or more individual processors (called cores in this sense) have been attached.
- In a many-core processor the number of cores is large enough that traditional multi-processor techniques are no longer efficient.
- 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, SIMD or multi-threading.
- Many applications do not realize yet large speedup factors: parallelizing algorithms and software is a major on-going research area.

	Memory	I/O
1_	TTT	
\langle	Netv	vork
V	S P P) S

Chip Multiprocessor (CMP)

-

< ロ > < 同 > < 回 > < 回 > :

Multi-core Architecture CPU Cache

CPU Cache (1/7)

Main Memory Index Data 0 xyz 1 pdq 2 abc 3 rgf

- A CPU cache is an auxiliary memory which is smaller, faster memory than the main memory and which stores copies of 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.



Multi-core Architecture CPU Cache

- 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).

	< □ >	▲ @ ▶ ▲ 문 ▶ ▲ 문 ▶ 문	4) Q (4		< □	1) 《 @) 《 문) 《 문) 문	£
(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	9 / 60	(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	10 / 60
٩	Multi-core Architecture CPU Cache			N	Aulti-core Architecture CPU Cache		
CPU Cache (3/7)				CPU Cache (4/7)			

CPU Cache (2/7)



- 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.



- 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

((Moreno	Maza)	

CS 4435 - CS 9624 11 / 60

3

500

(Moreno Maza)

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.



- 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).

(Moreno Maza)	Introduction t

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

15 / 60

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 Coherence (1/6)



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

16 / 60

Multi-core Archi	cecture CPU Coherence	Multi-core Architecture	CPU Coherence
Cache Coherence (2/6)		Cache Coherence (3/6)	



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



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

	< 🗆		🛯 🔊 < 🖓		4	ロト・日本・日本・日本・日本	ছ ৩৭৫
(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	17 / 60	(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	18 / 60
	Multi-core Architecture CPU Coherence				Multi-core Architecture CPU Coherence		
Cache Coherence	e (4/6)			Cache Coherence	(5/6)		



Figure: Processor P_2 issues a write x=5



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

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ―臣 – のへで

20 / 60

Multi-core Architecture CPU Coherence

Cache Coherence (6/6)



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

Multi-core Architecture CPU Coherence

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.

< ロ > < 邑 > < 臣 > < 臣 > < 臣 > < 臣 > < 臣 > のへの					< 🗆	▶ ★@ ▶ ★ 国 ▶ ★ 国 ▶ → 国	596
(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	21 / 60	(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	22 / 60
М	ulti-core Architecture CPU Coherence			N	Iulti-core Architecture CPU Coherence		

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

(日) (日) (日) (日) (日) (日) (日)

24 / 60

Introduction to Multicore Programming CS 4435 - CS 9624 (Moreno Maza)

23 / 60

Concurrency Platforms	Concurrency Platforms
an	Concurrency Platforms
Multi-core Architecture	 Programming directly on processor cores is painful and
• Multi-core processor	 Concurrency platforms
CPU CacheCPU Coherence	 abstract processor cores, handles synchronization, comr protocols

- 2 Concurrency Platforms
 - PThreads
 - TBB
 - Open MP
 - Cilk ++

(Mo

- Race Conditions and Cilkscreen
- MMM in Cilk++

error-prone.

- nunication protocols
- (optionally) perform load balancing
- Examples of concurrency platforms:
 - Pthreads
 - Threading Building Blocks (TBB)
 - OpenMP
 - Cilk++
- We use an implementation of the Fibonacci sequence $F_{n+2} = F_{n+1} + F_n$ to compare these four concurrency platforms.

	< □ >	 < □ > < □ > < □ > < □ 	• ৩৫৫		< □	 < □> < □> < □> Ξ 	596
(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	25 / 60	(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	26 / 60
	Concurrency Platforms				Concurrency Platforms PThreads		
Fibonacci Executi	ion			PThreads			



- Pthreads is a POSIX standard for threads, communicating though shared memory.
- Pthreads defines a set of C programming language types, functions and constants.
- It is implemented with a pthread.h header and a thread library.
- Programmers can use Pthreads to create, manipulate and manage threads.
- In particular, programmers can synchronize between threads using mutexes, condition variables and semaphores.
- This is a Do-it-yourself concurrency platform: programmers have to map threads onto the computer resources (static scheduling).

	< 🗆	▶ ▲@ ▶ ▲콜 ▶ ▲콜 ▶ _ 콜	$\mathcal{O}\mathcal{A}\mathcal{O}$		۹ 🗖	→ ▲圖 → ▲ 国 → ▲ 国 → 二日	🛯 ৩৫৫
reno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	27 / 60	(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	28 / 60

Plan

Concurrency Platforms PThreads

//object to set thread attributes (NULL for default)

//returned identifier for the new thread

//routine executed after creation

//a single argument passed to func

//identifier of thread to wait for

*WinAPI threads provide similar functionality.

//terminating thread's status (NULL to ignore)

Key PThread Function

const pthread_attr_t *attr,

void *(*func)(void *),

) //returns error status

pthread_t thread,

) //returns error status

int pthread_join (

void **status

int pthread_create(

void *arg

pthread_t *thread,

PThreads

- **Overhead:** The cost of creating a thread is more than 10,000 cycles. This enforces coarse-grained concurrency. (Thread pools can help.)
- Scalability: Fibonacci code gets about 1.5 speedup for 2 cores for computing fib(40).
 - Indeed the thread creation overhead is so large that only one thread is used, see below.
 - Consequently, one needs to rewrite the code for more than 2 cores.

Simplicity: Programmers must engage in error-prone protocols in order to schedule and load-balance.

	•		うへで			 → < @ > < E > < E > < 	€
(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	29 / 60	(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	30 / 60
	Concurrency Platforms PThreads				Concurrency Platforms PThreads		
				int main (int argc char	*=rgv[])		
<pre>#include <stdio.h></stdio.h></pre>				f	wargv [])		
<pre>#include <stdlib.h></stdlib.h></pre>				nthread t thread.			
<pre>#include <pthread.h></pthread.h></pre>				thread args args:			
				int status:			
int fib(int n)				int result:			
{				int thread result:			
if (n < 2) return n;				if (argc < 2) return	1:		
else {				int n = atoi(argv[1])	;		
int x = fib(n-1);				if (n < 30) result =	fib(n);		
int $y = fib(n-2);$				else {			
return x + y;				args.input = n-1;			
}				status = pthread_cr	eate(&thread,		
}					NULL,		
int input:					thread_func,		
int input;					<pre>(void*) &args);</pre>		
lift output,				// main can continu	e executing		
J thread_args,				<pre>result = fib(n-2);</pre>			
void *thread func (void	*ntr)			<pre>// Wait for the thr</pre>	ead to terminate.		
{	· ·poi)			pthread_join(thread	, NULL);		
int i = ((thread args))	*) ptr)->input:			result += args.outp	ut;		
((thread args *) ptr)-	>output = fib(i):			}			
return NULL;				printf("Fibonacci of	%d is %d.\n", n, result);		
}				return 0;			
				ł			
(Moreno Maza)	Introduction to Multicore Programming		31 / 60	(Moreno Maza)	Introduction to Multicore Programming		32 / 60

TBB (1/2)

Concurrency Platforms TBB

TBB (2/2)

- A C++ library that run on top of native threads
- Programmers specify tasks rather than threads:
 - Tasks are objects. Each task object has an input parameter and an output parameter.
 - One needs to define (at least) the methods: one for creating a task and one for executing it.
 - Tasks are launched by a spawn or a spawn_and_wait_for_all statement.
- Tasks are automatically load-balanced across the threads using the work stealing principle.
- TBB Developed by Intel and focus on performance

- $\bullet~\mathsf{TBB}$ provides many C++ templates to express common patterns simply, such as:
 - parallel_for for loop parallelism,
 - parallel_reduce for data aggregation
 - pipeline and filter for software pipelining
- TBB provides concurrent container classes which allow multiple threads to safely access and update items in the container concurrently.
- TBB also provides a variety of mutual-exclusion library functions, including locks.

	▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ ■ りへぐ		< □ >	 <	
(Moreno Maza) Introduction to Multicore Program Concurrency Platforms TBB	ming CS 4435 - CS 9624 33 / 60	(Moreno Maza)	Introduction to Multicore Programming Concurrency Platforms Open MP	CS 4435 - CS 9624 34 / 60	
<pre>(Moreno Maza) Introduction to Multicore Program Concurrency Platforms TBB class FibTask: public task { public: const long n; long* const sum; FibTask(long n_, long* sum_) : n(n_), sum(sum_) {} task* execute() { if(n < 2) { *sum = n; } else { long x, y; FibTask& a = *new(allocate_child())</pre>	iming CS 4435 - CS 9624 33 / 60	 Several compilers available, both open-source and Visual Studio. Runs on top of native threads Linguistic extensions to C/C++ or Fortran in the form of compiler pragmas (compiler directives): # pragma omp task shared(x) implies that the next statement is a independent task; 			
<pre>Fiblask(h=1,&x); FibTask& b = *new(allocate_child())</pre>		 moreover si other pragr data aggreg Supports loop parallelism with OpenMP provise mutual-exclusion 	haring of memory is managed explici- nas express directives for scheduling, gation. parallelism and, more recently in h dynamic scheduling. des a variety of synchronization c on locks, etc.)	tly loop parallelism amd Version 3.0, task onstructs (barriers,	

Concurrency Platforms Open MP	Concurrency Platforms Cilk ++			
	From Cilk to Cilk++			
<pre>int fib(int n) { if (n < 2) return n; int x, y; #pragma omp task shared(x) x = fib(n - 1); #pragma omp task shared(y) y = fib(n - 2); #pragma omp taskwait return x+y; }</pre>	 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 lead to Cilk++, developed by Cilk Arts, an MIT spin-off, which was acquired by Intel in July 2009. Today, it can be freely downloaded. The place where to start is http://www.cilk.com/ Cilk is still developed at MIT http://supertech.csail.mit.edu/cilk/ 			
(Moreno Maza) Introduction to Multicore Programming CS 4435 - CS 9624 37 / 60	(Moreno Maza) Introduction to Multicore Programming CS 4435 - CS 9624 38 / 60			
Cilk ++	Nested Parallelism in Cilk ++			

```
• 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.

```
int fib(int n)
ſ
```

```
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.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

(Moreno Maza)



The iterations of a cilk_for loop may execute in parallel.

- 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



Concurrency Platforms Cilk ++

Scheduling (2/3)

Concurrency Platforms Cilk ++

Scheduling (2/3)

- Cilk/Cilk++ randomized work-stealing scheduler load-balances the computation at run-time. Each processor maintains a ready deque:
 - A ready deque is a double ended queue, where each entry is a procedure instance that is ready to execute.
 - Adding a procedure instance to the bottom of the deque represents a procedure call being spawned.
 - A procedure instance being deleted from the bottom of the deque represents the processor beginning/resuming execution on that procedure.
 - Deletion from the top of the deque corresponds to that procedure instance being stolen.
- A mathematical proof guarantees near-perfect linear speed-up on applications with sufficient parallelism, as long as the architecture has sufficient memory bandwidth.
- A spawn/return in Cilk is over 100 times faster than a Pthread create/exit and less than 3 times slower than an ordinary C function call on a modern Intel processor.



▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ―臣 – のへで

(Moreno Maza)	Introducti Concurrency P	on to Multicore Program latforms Cilk ++	mming	CS 4435 - CS 9624	45 / 60	(Moreno Maza)	Introductio Concurrency Pla	on to Multicore Programm atforms Cilk ++	ing CS 4	435 - CS 9624	46 / 60
-						-	_	_			









- Iterations of a cilk_for should be independent.
- Between a cilk_spawn and the corresponding cilk_sync, the code of the spawned child should be independent of the code of the parent, including code executed by additional spawned or called children.
- The arguments to a spawned function are evaluated in the parent before the spawn occurs.



CS 4435 - CS 9624 55 / 60

(Moreno Maza)

в

Introduction to Multicore Programming

CS 4435 - CS 9624 56 / 60

Concurrency Platforms Race Conditions and Cilkscreen

Race Bugs (3/3)

• Watch out for races in packed data structures such as:

```
struct{
    char a;
```

```
char b;
```

(Moreno Maza)

```
}
```

Updating x.a and x.b in parallel can cause races.

- If an ostensibly deterministic Cilk++ program run on a given input could possibly behave any differently than its serialization, Cilkscreen race detector guarantees to report and localize the offending race.
- Employs a regression-test methodology (where the programmer provides test inputs) and dynamic instrumentation of binary code.
- Identifies files-names, lines and variables involved in the race.
- Runs about 20 times slower than real-time.

```
template<typename T> void multiply_iter_par(int ii, int jj, int kk,
        T* C)
{
        cilk_for(int i = 0; i < ii; ++i)
        for (int k = 0; k < kk; ++k)
            cilk_for(int j = 0; j < jj; ++j)
                 C[i * jj + j] += A[i * kk + k] + B[k * jj + j];
}
```

Does not scale up well due to a poor locality and uncontrolled granularity.

• Runs about 20 time				< □		-	
(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	57 / 60	(Moreno Maza)	Introduction to Multicore Programming	CS 4435 - CS 9624	58 / 60
Concurrency Platforms MMM in Cilk++			Co	oncurrency Platforms MMM in Cilk++			

}

CS 4435 - CS 9624

```
template<typename T> void multiply_rec_seq_helper(int i0, int i1, int j0,
    int j1, int k0, int k1, T* A, ptrdiff_t lda, T* B, ptrdiff_t ldb, T* C,
    ptrdiff_t ldc)
{
    int di = i1 - i0;
    int dj = j1 - j0;
    int dk = k1 - k0;
    if (di >= dj && di >= dk && di >= RECURSION_THRESHOLD) {
        int mi = i0 + di / 2;
        multiply_rec_seq_helper(i0, mi, j0, j1, k0, k1, A, lda, B, ldb, C, ldc);
        multiply_rec_seq_helper(mi, i1, j0, j1, k0, k1, A, lda, B, ldb, C, ldc);
   } else if (dj >= dk && dj >= RECURSION_THRESHOLD) {
        int mj = j0 + dj / 2;
        multiply_rec_seq_helper(i0, i1, j0, mj, k0, k1, A, lda, B, ldb, C, ldc);
        multiply_rec_seq_helper(i0, i1, mj, j1, k0, k1, A, lda, B, ldb, C, ldc);
    } else if (dk >= RECURSION_THRESHOLD) {
        int mk = k0 + dk / 2;
        multiply_rec_seq_helper(i0, i1, j0, j1, k0, mk, A, lda, B, ldb, C, ldc);
        multiply_rec_seq_helper(i0, i1, j0, j1, mk, k1, A, lda, B, ldb, C, ldc);
    } else {
        for (int i = i0; i < i1; ++i)
           for (int k = k0; k < k1; ++k)
                for (int j = j0; j < j1; ++j)
                    C[i * ldc + j] += A[i * lda + k] * B[k * ldb + j];
                                                     ▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶   臣 … のへで
```

Introduction to Multicore Programming

```
template<typename T> inline void multiply_rec_seq(int ii, int jj, i)
    T* B, T* C)
{
```

multiply_rec_seq_helper(0, ii, 0, jj, 0, kk, A, kk, B, jj, C, j)

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.

Introduction to Multicore Programming

#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

(Moreno Maza)