15-210: Parallelism in the Real World

- · Types of paralellism
- · Parallel Thinking
- · Nested Parallelism
- Examples (Cilk, OpenMP, Java Fork/Join)
- Concurrency

15-210 Page1

<u>Cray-1 (1976): the world's most</u> <u>expensive love seat</u>



15-210

<u>Data Center: Hundred's of</u> <u>thousands of computers</u>



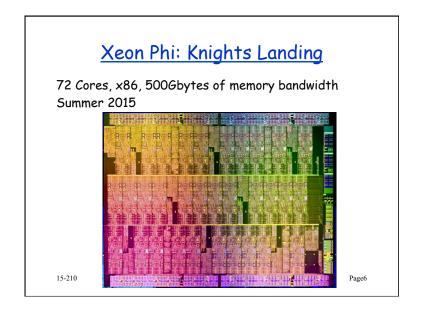


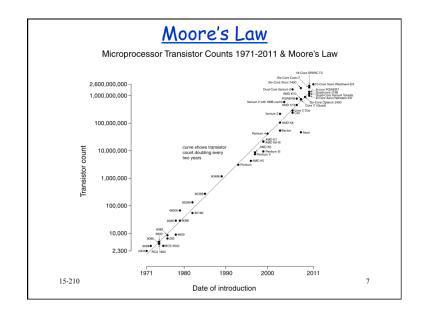


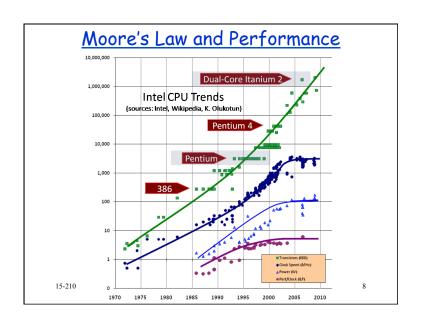
3

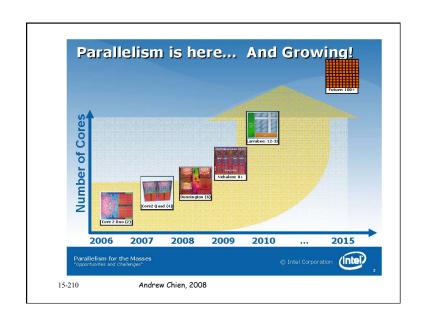


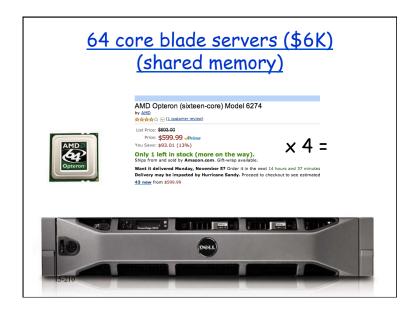
Telegrated Fabric Processor Package Integrated Fabric Processor Pac















Intel Has a 48-Core Chip for Smartphones and Tablets

By Wolfgang Gruener OCTOBER 31, 2012 9:20 AM - Source: Computerworld

Intel has developed a prototype of a 48-core processor for smartphones. Before you ask: No, you can't buy a 48-core smartphone next year.



15-210

Key Challenge: Software (How to Write Parallel Code?)

At a high-level, it is a two step process:

- Design a work-efficient, low-span parallel algorithm
- Implement it on the target hardware

In reality: each system required different code because programming systems are immature

- Huge effort to generate efficient parallel code.
 - Example: Quicksort in MPI is <u>1700 lines</u> of code, and about the same in CUDA
- Implement one parallel algorithm: a whole thesis. Take 15-418 (Parallel Computer Architecture and Prog.)

15-210

Parallel Hardware

Many forms of parallelism

- Supercomputers: large scale, shared memory
- Clusters and data centers: large-scale, distributed memory
- Multicores: tightly coupled, smaller scale
- GPUs, on chip vector units
- Instruction-level parallelism

Parallelism is important in the real world.

15-210

15-210 Approach

Enable parallel thinking by raising abstraction level

- I. Parallel thinking: Applicable to many machine models and programming languages
- II. Reason about correctness and efficiency of algorithms and data structures.

15-210 16

Parallel Thinking

Recognizing true dependences: unteach sequential programming.

Parallel algorithm-design techniques

- Operations on aggregates: map/reduce/scan
- Divide & conquer, contraction
- Viewing computation as DAG (based on dependences)

Cost model based on work and span

15-210

Quicksort from Aho-Hopcroft-Ullman (1974)

procedure QUICKSORT(S):

if S contains at most one element then return S else

begin

choose an element **a** randomly from **S**; **let S**₁, **S**₂ and **S**₃ be the sequences of elements in **S** less than, equal to, and greater than **a**, respectively; **return** (QUICKSORT(**S**₁) followed by **S**₂ followed by QUICKSORT(**S**₃))

end

15-210



Page 18

Quicksort from Sedgewick (2003)

```
public void quickSort(int[] a, int left, int right) {
    int i = left-1; int j = right;
    if (right <= left) return;</pre>
    while (true) {
      while (a[++i] < a[right]);
      while (a[right] < a[--j])
                                         Algorithms
        if (j==left) break;
      if (i >= j) break;
                                              INJava
      swap(a,i,j); }
    swap(a, i, right);
    quickSort(a, left, i - 1);
    quickSort(a, i+1, right); }
  15-210
                                                  Page 19
```

Styles of Parallel Programming

Data parallelism/Bulk Synchronous/SPMD Nested parallelism: what we covered Message passing Futures (other pipelined parallelism) General Concurrency

15-210

Page20

Nested Parallelism

Nested Parallelism =

arbitrary nesting of parallel loops + fork-join

- Assumes no synchronization among parallel tasks except at joint points.
- Deterministic if no race conditions

Advantages:

- Good schedulers are known
- Easy to understand, debug, and analyze costs
- Purely functional, or imperative...either works

15-210 Page21

Serial Parallel DAGs

Dependence graphs of nested parallel computations are series parallel



Two tasks are parallel if not reachable from each other. A data race occurs if two parallel tasks are involved in a race if they access the same location and at least one is a write.

15-210 Page22

Nested Parallelism: parallel loops

Nested Parallelism: fork-join

Nested Parallelism: fork-join

Cilk vs. what we've covered

```
ML:
               val(a,b) = par(fn() => f(x), fn() => g(y))
Psuedocode: val (a,b) = (f(x) || g(y))
                cilk\_spawn a = f(x);
Cilk:
                                                  Fork Join
                b = g(y);
                cilk_sync;
                S = map f A
ML:
                                                   Map
Psuedocode: S = \langle f x : x \text{ in } A \rangle
                cilk_for (int i = 0; i < n; i++)
Cilk:
                  S[i] = f(A[i])
15-210
                                                        Page26
```

Cilk vs. what we've covered

Page25

15-210 Page27

Example Cilk

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

Example OpenMP: Numerical Integration Mathematically, we know that: We can approximate the $F(x) = 4.0/(1+x^2)$ integral as a sum of rectangles: $\sum_{i=0}^{N} F(x_i) \Delta x \approx \pi$ where each rectangle has width Δx and height $F(x_i)$ at Χ 15-210 the middle of interval i.

The C/openMP code for Approx. PI #include <omp.h> static long num steps = 100000; double step: void main () Private clause { int i; double x, pi, sum = 0.0; creates data local to step = 1.0/(double) num steps; a thread #pragma omp parallel for private(i, x) reduction(+:sum) for $(i=0;i\leq num steps; i++)$ x = (i+0.5)*step: sum = sum + 4.0/(1.0+x*x);pi = step * sum; Reduction used to manage dependencies

The C code for Approximating PI static long num steps = 100000; double step; void main () int i; double x, pi, sum = 0.0; step = 1.0/(double) num_steps; x = 0.5 * step;for $(i=0;i\leq num steps; i++)$ x+=step; sum += 4.0/(1.0+x*x);

pi = step * sum;

```
Example: Java Fork/Join
class Fib extends FJTask {
  volatile int result; // serves as arg and result
  Fib(int _n) { n = _n; }
  public void run() {
     if (n \le 1) result = n;
      else if (n <= sequentialThreshold) number = seqFib(n);</pre>
        Fib f1 = new Fib(n - 1);
        Fib f2 = new Fib(n - 2);
        coInvoke(f1, f2);
        result = f1.result + f2.result;
}
15-210
                                                        Page32
```

Cost Model (General)

Compositional:

Work: total number of operations

- costs are added across parallel calls

Span: depth/critical path of the computation

- Maximum span is taken across forked calls

Parallelism = Work/Span

- Approximately # of processors that can be effectively used.

15-210 Page33

Combining costs (Nested Parallelism)

Combining for parallel for:

$$W_{\text{pexp}}(\text{pfor ...}) = \sum_{i=0}^{n-1} W_{\text{exp}}(f(i))$$
 work

$$D_{\text{pexp}}(\text{pfor }...) = \max_{i=0}^{n-1} D_{\text{exp}}(f(i))$$
 span

15-210 34

Why Work and Span

Simple measures that give us a good sense of efficiency (work) and scalability (span).

Can schedule in O(W/P + D) time on P processors.

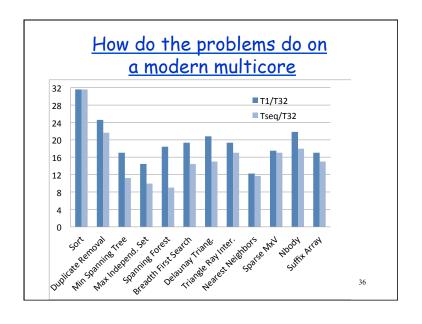
This is within a constant factor of optimal.

Goals in designing an algorithm

- Work should be about the same as the sequential running time. When it matches asymptotically we say it is work efficient.
- 2. Parallelism (W/D) should be polynomial. $O(n^{1/2})$ is probably good enough

15-210

35



Styles of Parallel Programming

Data parallelism/Bulk Synchronous/SPMD Nested parallelism: what we covered Message passing Futures (other pipelined parallelism)

General Concurrency

15-210 Page37

Parallelism vs. Concurrency

- Parallelism: using multiple processors/cores running at the same time. Property of the machine
- Concurrency: non-determinacy due to interleaving threads. Property of the application.

		Concurrency	
		sequential	concurrent
Parallelism	serial	Traditional programming	Traditional OS
	parallel	Deterministic parallelism	General parallelism

15-210 38

Concurrency: Stack Example 1

```
struct link {int v; link* next;}
struct stack {
  link* headPtr;
  void push(link* a) {
    a->next = headPtr;
    headPtr = a; }
  link* pop() {
    link* h = headPtr;
    if (headPtr != NULL)
        headPtr = headPtr->next;
    return h;}
}
```

Concurrency: Stack Example 1

```
struct link {int v; link* next;}
struct stack {
   link* headPtr;
   void push(link* a) {
      a->next = headPtr;
      headPtr = a; }
   link* pop() {
      link* h = headPtr;
      if (headPtr != NULL)
        headPtr = headPtr->next;
      return h;}
}
```

concurrency: Stack Example 1 struct link {int v; link* next;} struct stack { link* headPtr; void push(link* a) { a->next = headPtr; headPtr = a; } link* pop() { link* h = headPtr; if (headPtr != NULL) headPtr = headPtr->next; return h;} }

concurrency: Stack Example 1 struct link {int v; link* next;} struct stack { link* headPtr; void push(link* a) { a->next = headPtr; headPtr = a; } link* pop() { link* h = headPtr; if (headPtr != NULL) headPtr = headPtr->next; return h;} }

```
bool CAS(ptr* addr, ptr a, ptr b) {
  atomic {
    if (*addr == a) {
        *addr = b;
        return 1;
    } else
        return 0;
}

A built in instruction on most processors:

CMPXCHG8B - 8 byte version for x86

CMPXCHG16B - 16 byte version

15-210

Page43
```

```
Concurrency: Stack Example 2
struct stack {
 link* headPtr;
 void push(link* a) {
     link* h = headPtr:
     a->next = h:
   while (!CAS(&headPtr, h, a)); }
 link* pop() {
   do {
     link* h = headPtr;
     if (h == NULL) return NULL;
     link* nxt = h->next;
   while (!CAS(&headPtr, h, nxt))}
   return h;}
 15-210
                                                44
```

concurrency: Stack Example 2 struct stack { link* headPtr; void push(link* a) { do { link* h = headPtr; a->next = h; while (!CAS(&headPtr, h, a)); } link* pop() { do { link* h = headPtr; if (h == NULL) return NULL; link* nxt = h->next; while (!CAS(&headPtr, h, nxt))} return h; } 15-210

```
Concurrency: Stack Example 2
struct stack {
  link* headPtr;
  void push(link* a) {
    do {
      link* h = headPtr:
      a \rightarrow next = h:
while (!CAS(&headPtr, h, a)); }
  link* pop() {
    do {
      link* h = headPtr;
      if (h == NULL) return NULL;
      link* nxt = h->next;
    while (!CAS(&headPtr, h, nxt))}
    return h;}
  15-210
                                                  47
```

concurrency: Stack Example 2 struct stack { link* headPtr; void push(link* a) { do { link* h = headPtr; a->next = h; while (!CAS(&headPtr, h, a)); } link* pop() { do { link* h = headPtr; if (h == NULL) return NULL; link* nxt = h->next; while (!CAS(&headPtr, h, nxt)) } return h; } }

```
Concurrency: Stack Example 2
struct stack {
  link* headPtr;
  void push(link* a) {
    do {
      link* h = headPtr:
      a->next = h:
while (!CAS(&headPtr, h, a)); }
  link* pop() {
    do {
      link* h = headPtr;
      if (h == NULL) return NULL;
      link* nxt = h->next;
    while (!CAS(&headPtr, h, nxt))}
    return h;}
  15-210
                                                 48
```

Concurrency: Stack Example 2' P1: x = s.pop(); y = s.pop(); s.push(x); P2: z = s.pop(); Before: A B C P2: h = headPtr; P2: nxt = h->next; P1: everything The ABA problem P2: CAS (&headPtr,h,nxt) Can be fixed with counter and 2CAS, but... 15-210 49

Concurrency: Stack Example 3

```
struct link {int v; link* next;}
struct stack {
    link* headPtr;
    void push(link* a) {
        atomic {
            a->next = headPtr;
            headPtr = a;      }}
link* pop() {
        atomic {
            link* h = headPtr;
            if (headPtr != NULL)
                 headPtr = headPtr->next;
                return h;}
}
```

Concurrency: Stack Example 3'

```
void swapTop(stack s) {
   link* x = s.pop();
   link* y = s.pop();
   push(x);
   push(y);
}
```

Queues are trickier than stacks.

15-210

51

Styles of Parallel Programming

Data parallelism/Bulk Synchronous/SPMD Nested parallelism: what we covered Message passing

Futures (other pipelined parallelism)
General Concurrency

15-210 Page52

Futures: Example

```
fun quickSort S = let
  fun qs([], rest) = rest
    | qs(h::T, rest) =
    let
      val L1 = filter (fn b => b < a) T
      val L2 = filter (fn b => b >= a) T
      in
        qs(L1,(a::(qs(L2, rest)))
      end

fun filter(f,[]) = []
    | filter(f,h::T) =
        if f(h) then(h::filter(f,T))
      else filter(f,T)

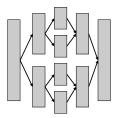
in
    qs(S,[])
end

Page53
```

Futures: Example

Quicksort: Nested Parallel

Parallel Partition and Append



Work = $O(n \log n)$

55

Span = $O(lg^2 n)$

15-210

Work = O(n log n)

Span = O(n)

15-210

56

Quicksort: Futures

Styles of Parallel Programming

Data parallelism/Bulk Synchronous/SPMD

- * Nested parallelism : what we covered Message passing
- * Futures (other pipelined parallelism)
- * General Concurrency

15-210 Page57

Question

How do we get nested parallel programs to work well on a fixed set of processors? Programs say nothing about processors.

Answer: good schedulers

15-210 Page59

Xeon 7500 Memory: up to 1 TB 4 of these 24 MB 24 MB L3 8 of these 8 of these **L2** 128 KB 128 KB 128 KB 128 KB 32 KB 32 KB 32 KB 32 KB L1 (P)15-210 Page58

Greedy Schedules

"Speedup versus Efficiency in Parallel Systems", Eager, Zahorjan and Lazowska, 1989

For any greedy schedule:

Parallel Time =
$$T_P \le \frac{W}{P} + D$$

15-210 60

Types of Schedulers

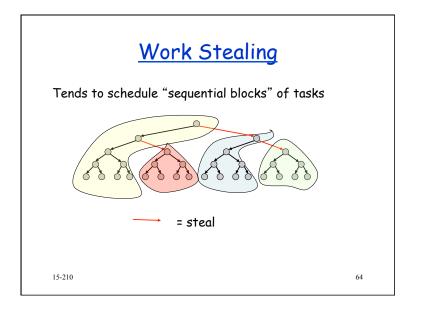
Bread-First Schedulers Work-stealing Schedulers Depth-first Schedulers

15-210

Page61

Breadth First Schedules Most naïve schedule. Used by most implementations of P-threads. O(n³) tasks Bad space usage, bad locality 15-210 62

Work Stealing P₁ P₂ P₃ P₄ old work queues P₁ P₂ P₃ P₄ old work queues push new jobs on "new" end pop jobs from "new" end If processor runs out of work, then "steal" from another "old" end Each processor tends to execute a sequential part of the computation. 15-210



Work Stealing Theory

For strict computations
Blumofe and Leiserson, 1999
of steals = O(PD)Space = $O(PS_1)$ S_1 is the sequential space
Acar, Blelloch and Blumofe, 2003
of cache misses on distributed caches
 M_1 + O(CPD) M_1 = sequential misses, C = cache size

15-210 65

Work Stealing Practice

Used in Cilk Scheduler

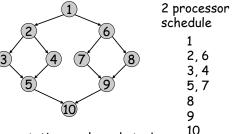
- Small overheads because common case of pushing/popping from local queue can be made fast (with good data structures and compiler help).
- No contention on a global queue
- Has good distributed cache behavior
- Can indeed require O(S1P) memory

Used in X10 scheduler, and others

15-210 66

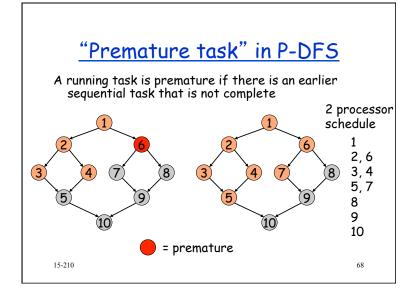
Parallel Depth First Schedules (P-DFS)

List scheduling based on Depth-First ordering



67

For strict computations a shared stack implements a P-DFS $\,$



P-DFS Theory

Blelloch, Gibbons, Matias, 1999
For any computation:
Premature nodes at any time = O(PD)Space = $S_1 + O(PD)$ Blelloch and Gibbons, 2004
With a shared cache of size $C_1 + O(PD)$ we have $M_p = M_1$

15-210 69

Conclusions

- lots of parallel languages
- lots of parallel programming styles
- high-level ideas cross between languages and styles
- scheduling is important

15-210 Page71

P-DFS Practice

Experimentally uses less memory than work stealing and performs better on a shared cache.

Requires some "coarsening" to reduce overheads

15-210 70