# Topic 14:     Parallelism

COS 320

Compiling Techniques

Princeton University
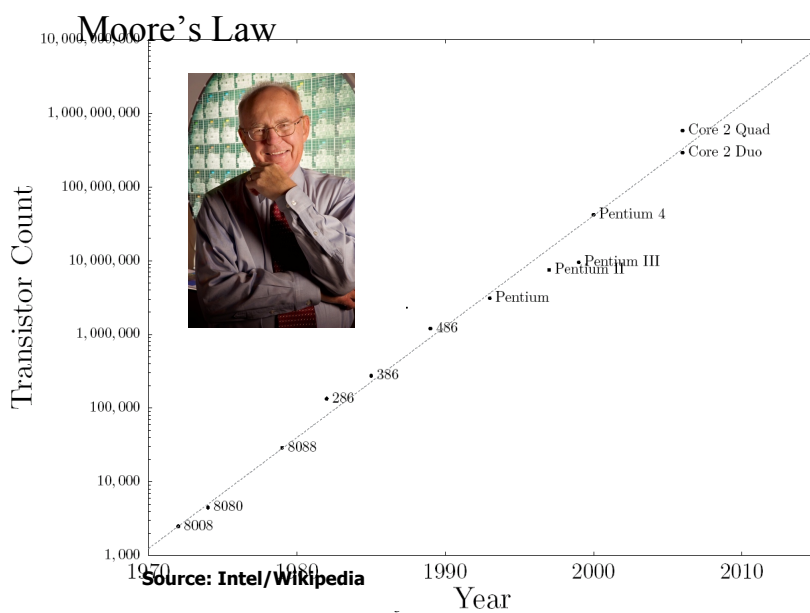Spring 2015

Prof. David August
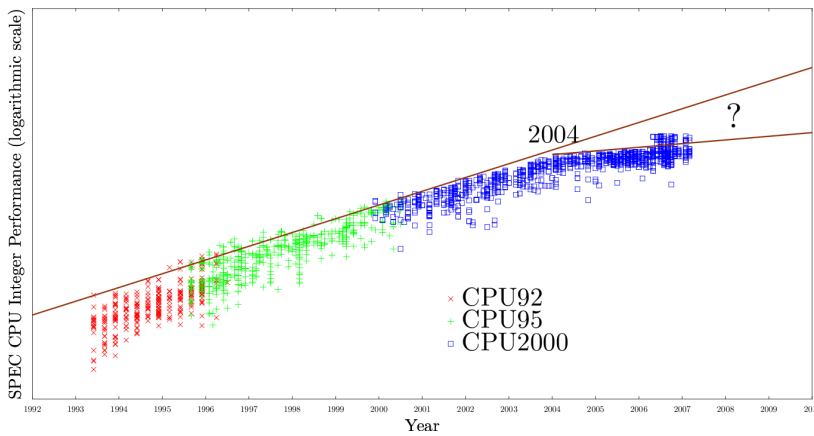
1

# Final Exam!

- Friday May 22 at 1:30PM in FRIEND 006
- Closed book
- One Front/Back 8.5x11

2

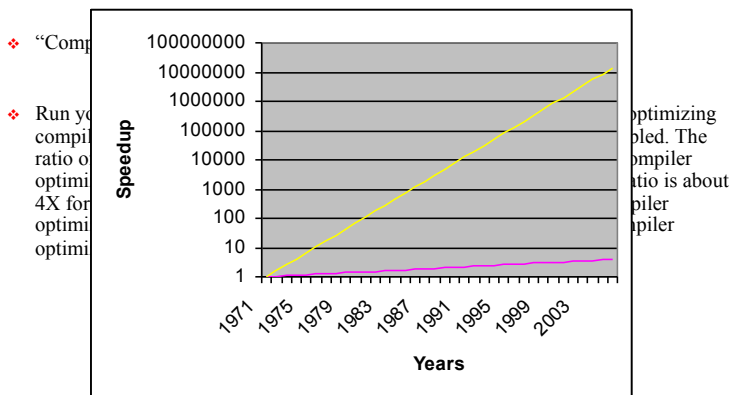# Moore's Law

## Single-Threaded Performance Not Improving

## What about Parallel Programming? –or-
## What is Good About the Sequential Model?

- ❖ Sequential is easier
  - » People think about programs sequentially
  - » Simpler to write a sequential program
- ❖ Deterministic execution
  - » Reproducing errors for debugging
  - » Testing for correctness
- ❖ No concurrency bugs
  - » Deadlock, livelock, atomicity violations
  - » Locks are not composable
- ❖ Performance extraction
  - » Sequential programs are portable
    - Ÿ Are parallel programs?  Ask GPU developers ☺
  - » Performance debugging of sequential programs straight-forward
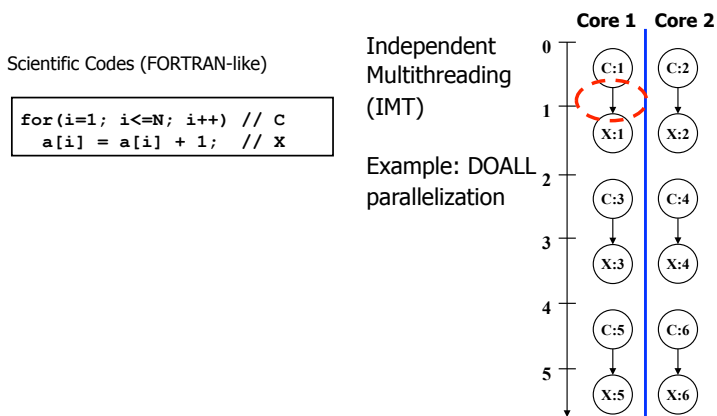
## Compilers are the Answer? - Proebsting's Law

- ❖ "Comp

- ❖ Run yo                                                       optimizing
  compil                                                        bled. The
  ratio o                                                       ompiler
  optimi                                                        atio is about
  4X for                                                        piler
  optimi                                                        piler
  optimi



**Conclusion – Compilers not about performance!**

# Are We Doomed?

## A Step Back in Time: Old Skool Parallelization

---

## Parallelizing Loops In Scientific Applications

Scientific Codes (FORTRAN-like)

```
for(i=1; i<=N; i++) // C
  a[i] = a[i] + 1;  // X
```

Independent Multithreading (IMT)

Example: DOALL parallelization

**Core 1   Core 2**

---

## What Information is Needed to Parallelize?

- ❖ Dependences within iterations are fine
- ❖ Identify the presence of cross-iteration data-dependences
  - » Traditional analysis is inadequate for parallelization. For instance, it does not distinguish between different executions of the same statement in a loop.
- ❖ Array dependence analysis enables optimization for parallelism in programs involving arrays.
  - » Determine pairs of iterations where there is a data dependence
  - » Want to know all dependences, not just yes/no

```
for(i=1; i<=N; i++) // C
  a[i] = a[i] + 1;  // X
```

```
for(i=1; i<=N; i++) // C
  a[i] = a[i-1] + 1;  // X
```

## Affine/Linear Functions

❖ f( $i_1, i_{2, \ldots,} i_n$) is <u>affine</u>, if it can be expressed as a sum of a constant, plus constant multiples of the variables. i.e.

$$f = c_0 + \sum_{i=1}^{n} c_i x_i$$

❖ Array subscript expressions are usually affine functions involving loop induction variables.

❖ Examples:
- » a[ i ]  affine
- » a[ i+j -1 ]  affine
- » a[ i*j ]  non-linear, not affine
- » a[ 2*i+1, i*j ]  linear/non-linear, not affine
- » a[ b[i] + 1 ]  non linear (indexed subscript), not affine

## Array Dependence Analysis

for (i = 1; i < 10; i++) {

   X[i] = X[i-1]

}

To find all the data dependences, we check if

1. X[i-1]  and X[i] refer to the same location;
2. different instances of X[i] refer to the same location.
- » For 1, we solve for i and i' in
  $1 \le i \le 10$, $1 \le i' \le 10$ and $i - 1 = i'$
- » For 2, we solve for i and i' in
  $1 \le i \le 10$, $1 \le i' \le 10$, $i = i'$ and $i \ne i'$ (between different dynamic accesses)

There is a dependence since there exist integer solutions to 1. e.g. (i=2, i'=1), (i=3,i'=2). 9 solutions exist.

There is no dependences among different instances of X[i] because 2 has no solutions!

## Array Dependence Analysis - Summary

❖ Array data dependence basically requires finding integer solutions to a system (often refers to as dependence system) consisting of equalities and inequalities.

❖ Equalities are derived from array accesses.

❖ Inequalities from the loop bounds.

❖ It is an integer linear programming problem.

❖ ILP is an NP-Complete problem.

❖ Several Heuristics have been developed.
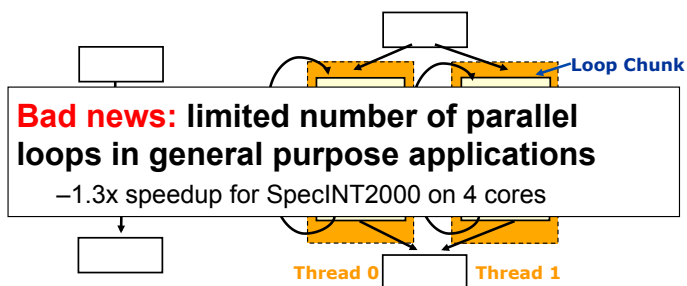- » Omega – U. Maryland

## Loop Parallelization Using Affine Analysis Is Proven Technology

- ❖ DOALL Loop
  - » No loop carried dependences for a particular nest
  - » Loop interchange to move parallel loops to outer scopes
- ❖ Other forms of parallelization possible
  - » DOAcross, DOpipe
- ❖ Optimizing for the memory hierarchy
  - » Tiling, skewing, etc.
- ❖ Real compilers available – KAP, Portland Group, gcc
- ❖ For better information, see
  - » http://gcc.gnu.org/wiki/Graphite?action=AttachFile&do=get&target=graphite_lambda_tutorial.pdf

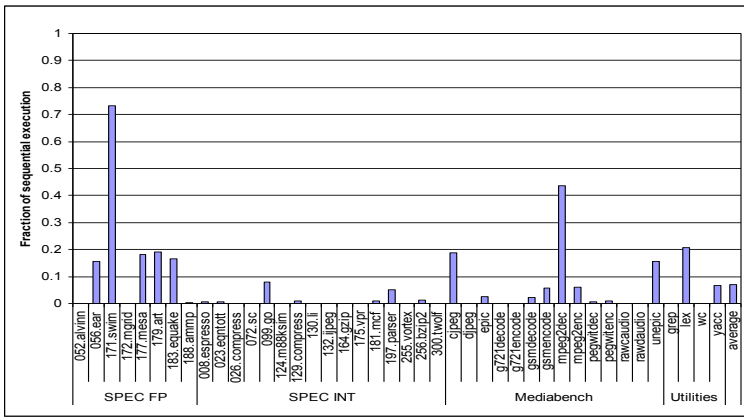# Back to the Present – Parallelizing C and C++ Programs

# Loop Level Parallelization

**Loop Chunk**
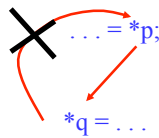
**Bad news:** limited number of parallel loops in general purpose applications

–1.3x speedup for SpecINT2000 on 4 cores

**Thread 0**    **Thread 1**

# DOALL Loop Coverage



Fraction of sequential execution

SPEC FP | SPEC INT | Mediabench | Utilities

# What's the Problem?

1. Memory dependence analysis

for (i=0; i<100; i++) {

. . . = *p;

*q = . . .

}

Memory dependence profiling
and speculative parallelization

# DOALL Coverage – Provable and Profiled



Fraction of sequential execution

☐ Profiled DOALL
■ Provable DOALL

**Still not good enough!**

SPEC | Utilities

## What's the Next Problem?

2. Data dependences

```
while (ptr != NULL) {

    . . .

    ptr = ptr->next;

    sum = sum + foo;
}
```
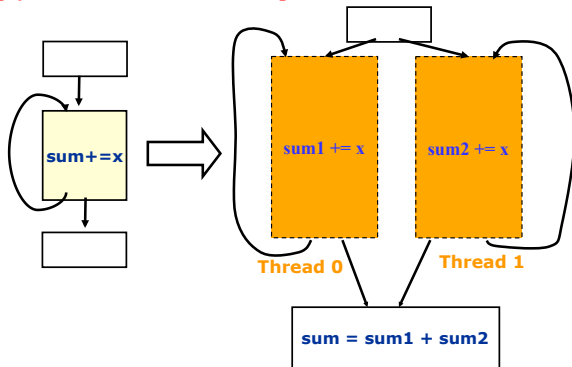
➡ Compiler transformations

## We Know How to Break Some of These Dependences – Recall ILP Optimizations

Apply accumulator variable expansion!



```
sum+=x     →     sum1 += x        sum2 += x
                 Thread 0         Thread 1

              sum = sum1 + sum2
```

## Data Dependences Inhibit Parallelization

- ❖ Accumulator, induction, and min/max expansion only capture a small set of dependences
- ❖ 2 options
  - » 1) Break more dependences – New transformations
  - » 2) Parallelize in the presence of dependences – more than DOALL parallelization
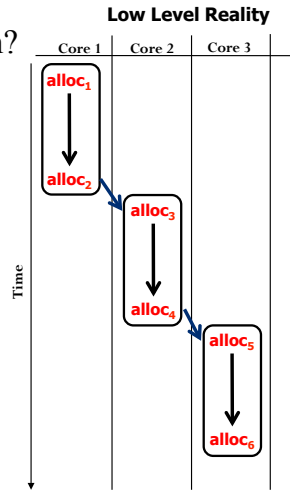- ❖ We will talk about both, but for now ignore this issue

## What's the Next Problem?

**3. C/C++ too restrictive**

```
char *memory;

void * alloc(int size);

void * alloc(int size) {
  void * ptr = memory;
  memory = memory + size;
  return ptr;
}
```
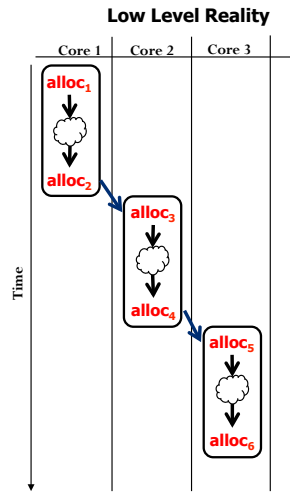


- 22 -

---

```
char *memory;

void * alloc(int size);

void * alloc(int size) {
  void * ptr = memory;
  memory = memory + size;
  return ptr;
}
```
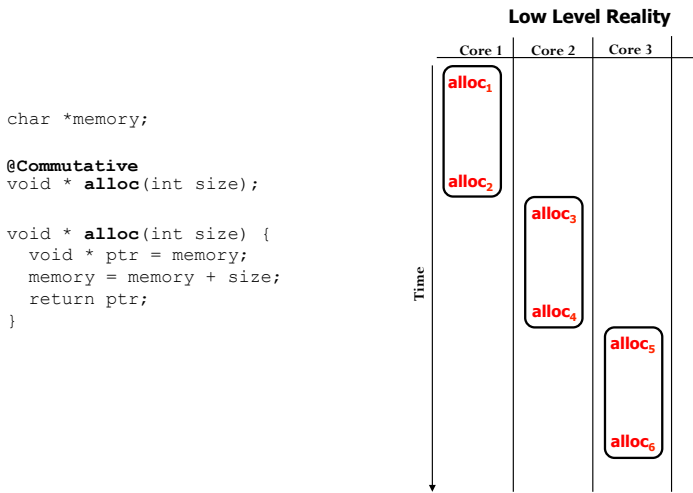
Loops cannot be parallelized even if computation is independent



- 23 -

---

## Commutative Extension

- ❖ Interchangeable call sites
  - » Programmer doesn't care about the order that a particular function is called
  - » Multiple different orders are all defined as correct
  - » Impossible to express in C
- ❖ Prime example is memory allocation routine
  - » Programmer does not care which address is returned on each call, just that the proper space is provided
- ❖ Enables compiler to break dependences that flow from 1 invocation to next forcing sequential behavior

- 24 -

**Low Level Reality**



```
char *memory;

@Commutative
void * alloc(int size);

void * alloc(int size) {
  void * ptr = memory;
  memory = memory + size;
  return ptr;
}
```

**Low Level Reality**



```
char *memory;

@Commutative
void * alloc(int size);

void * alloc(int size) {
  void * ptr = memory;
  memory = memory + size;
  return ptr;
}
```

Implementation dependences should not cause serialization.

# What is the Next Problem?
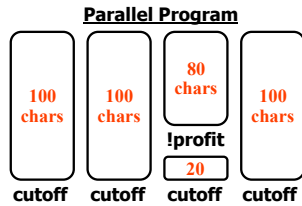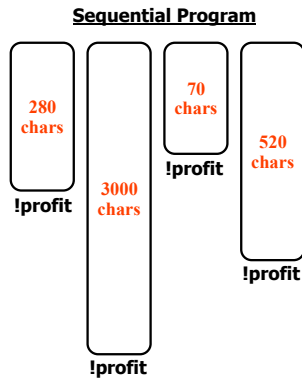
❖ 4. **C does not allow any prescribed non-determinism**
   » Thus sequential semantics must be assumed even though they not necessary
   » Restricts parallelism (useless dependences)
❖ Non-deterministic branch → programmer does not care about individual outcomes
   » They attach a probability to control how statistically often the branch should take
   » Allow compiler to tradeoff 'quality' (e.g., compression rates) for performance
      Ÿ When to create a new dictionary in a compression scheme

```
#define CUTOFF 100
dict = create_dict();
count = 0;
while((char = read(1))) {
  profitable =
        compress(char, dict)
  if (!profitable) {
    dict=restart(dict);
  } if (count == CUTOFF){
    dict=restart(dict);
    count=0;
  }

  count++;
}
finish_dict(dict);
```



**Sequential Program**

280 chars !profit
3000 chars !profit
70 chars !profit
520 chars !profit

**Parallel Program**

100 chars cutoff
100 chars cutoff
80 chars !profit 20 cutoff
100 chars cutoff

- 28 -

**2-Core Parallel Program**

```
dict = create_dict();
while((char = read(1))) {
  profitable =
        compress(char, dict)

  @YBRANCH(probability=.01)
  if (!profitable) {
    dict = restart(dict);
  }
}
finish_dict(dict);
```



1500 chars !prof / 1000 chars cutoff
2500 chars cutoff
800 chars !prof / 1700 chars cutoff
1200 chars !prof / 1300 chars cutoff

**Reset every 2500 characters**

**64-Core Parallel Program**

Compilers are best situated to make the tradeoff between output quality and performance

100 chars cutoff
100 chars cutoff
80 chars !prof 20 cutoff
100 chars cutoff

**Reset every 100 characters**

- 29 -

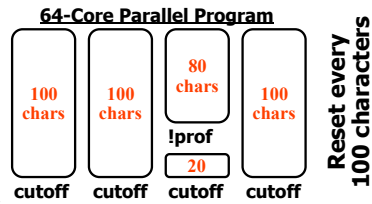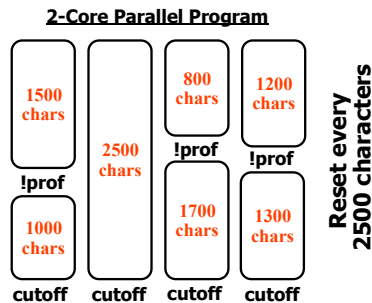## Capturing Output/Performance Tradeoff: *Y-Branches in 164.gzip*

```
dict = create_dict();
while((char = read(1))) {
  profitable =
        compress(char, dict)

  @YBRANCH(probability=.00001)
  if (!profitable) {
    dict = restart(dict);
  }
}
finish_dict(dict);
```
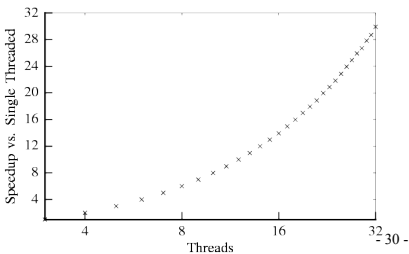
```
#define CUTOFF 100000
dict = create_dict();
count = 0;
while((char = read(1))) {
  profitable =
        compress(char, dict)

  if (!profitable)
    dict=restart(dict);
  if (count == CUTOFF){
    dict=restart(dict);
    count=0;
  }

  count++;
}
finish_dict(dict);
```



Speedup vs. Single Threaded vs. Threads

- 30 -

# 256.bzip2

```
unsigned char *block;
int last_written;
                                        doReversibleTransform() {
compressStream(in, out) {                 .
  while (True) {                          sortIt();
    loadAndRLEsource(in);                 ..
    if (!last) break;                    }

    doReversibleTransform();
                                        sortIt() {
    sendMTFValues(out);                   .
  }                                       printf(...);
}                                         ..
                                        }
```

Parallelization techniques must look inside function calls to
expose operations that cause synchronization.

# 197.parser

```
batch_process() {                       char *memory;
  while(True) {
    sentence = read();                  void * xalloc(int size) {
    if (!sentence) break;                 void * ptr = memory;
                                          memory = memory + size;
    parse(sentence);                      return ptr;
                                        }
    print(sentence);

  }
}
```

**High-Level View:**

Parsing a sentence is
independent of any other
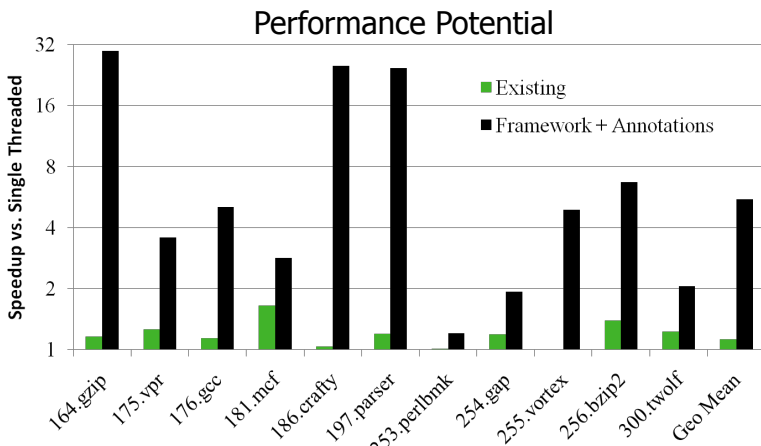sentence.

**Low-Level Reality:**

Implementation dependences
inside functions called by
*parse* lead to large
sequential regions.

| | LoC Changed | Increased Scope | Commutative | Y-Branch | Nested Parallel | Iter. Inv. Value Spec. | Loop Alias Spec. | Programmer Mod. |
|---|---|---|---|---|---|---|---|---|
| 164.gzip | 26 | x | | x | | | | x |
| 175.vpr | 1 | | x | | | x | x | |
| 176.gcc | 18 | x | x | | | | x | x |
| 181.mcf | 0 | | | | x | | | |
| 186.crafty | 9 | x | x | | x | x | x | |
| 197.parser | 3 | x | x | | | | | |
| 253.perlbmk | 0 | x | | | | x | x | |
| 254.gap | 3 | x | x | | | | x | |
| 255.vortex | 0 | x | | | | x | x | |
| 256.bzip2 | 0 | x | | | | | x | |
| 300.twolf | 1 | x | x | | | | x | |

**Modified only 60 LOC out of ~500,000 LOC**

## Performance Potential



Chart: Speedup vs. Single Threaded (y-axis: 1, 2, 4, 8, 16, 32) for benchmarks 164.gzip, 175.vpr, 176.gcc, 181.mcf, 186.crafty, 197.parser, 253.perlbmk, 254.gap, 255.vortex, 256.bzip2, 300.twolf, Geo Mean. Legend: ■ Existing, ■ Framework + Annotations

**What prevents the automatic extraction of parallelism?**

~~Lack of an Aggressive Compilation Framework~~
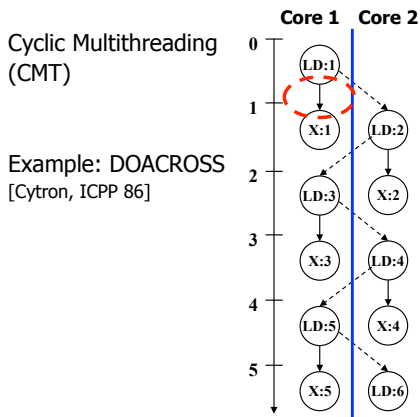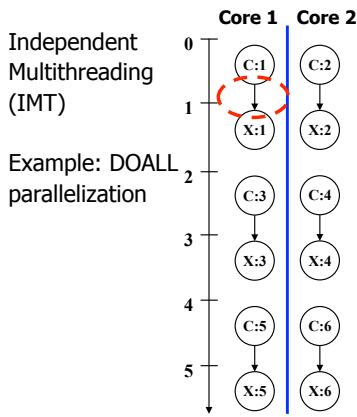
~~Sequential Programming Model~~

---

# What About Non-Scientific Codes???

Scientific Codes (FORTRAN-like)

```
for(i=1; i<=N; i++) // C
   a[i] = a[i] + 1;  // X
```

General-purpose Codes (legacy C/C++)
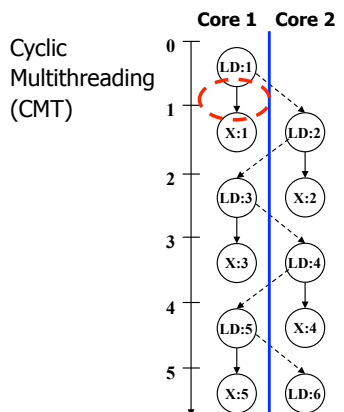
```
while(ptr = ptr->next)     // LD
   ptr->val = ptr->val + 1; // X
```
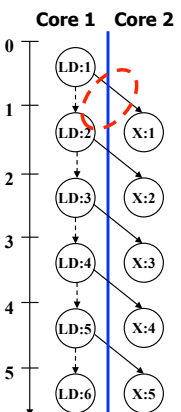
Independent Multithreading (IMT)

Example: DOALL parallelization

Cyclic Multithreading (CMT)

Example: DOACROSS [Cytron, ICPP 86]

---

# Alternative Parallelization Approaches

```
while(ptr = ptr->next)     // LD
   ptr->val = ptr->val + 1; // X
```
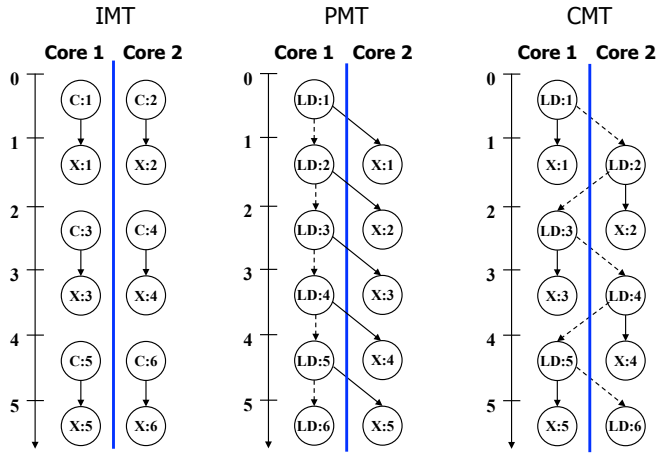
Cyclic Multithreading (CMT)

Pipelined Multithreading (PMT)

Example: DSWP **[PACT 2004]**

# Comparison: IMT, PMT, CMT

# Comparison: IMT, PMT, CMT



lat(comm) = 1:   1 iter/cycle          1 iter/cycle              1 iter/cycle

# Comparison: IMT, PMT, CMT



lat(comm) = 1:   1 iter/cycle          1 iter/cycle              1 iter/cycle
lat(comm) = 2:   1 iter/cycle          1 iter/cycle              0.5 iter/cycle

## Comparison: IMT, PMT, CMT

**Thread-local Recurrences ➔ Fast Execution**

| | IMT | | PMT | | CMT | |
|---|---|---|---|---|---|---|
| | **Core 1** | **Core 2** | **Core 1** | **Core 2** | **Core 1** | **Core 2** |

IMT:
Core 1: C:1, X:1, C:3, X:3, C:5, X:5
Core 2: C:2, X:2, C:4, X:4, C:6, X:6

PMT:
Core 1: LD:1, LD:2, LD:3, LD:4, LD:5, LD:6
Core 2: X:1, X:2, X:3, X:4

CMT:
Core 1: LD:1, X:1, LD:3, X:3
Core 2: LD:2, X:2

**Cross-thread Dependences ➔ Wide Applicability**

- 40 -

## Our Objective: Automatic Extraction of Pipeline Parallelism using DSWP

197.parser

Speedup vs. Single Threaded — Threads (1, 2, 4, 8, 16, 32, 64)

**Find English Sentences** → **Parse Sentences (95%)** → **Emit Results**

**Decoupled Software Pipelining**

**PS-DSWP (Spec DOALL Middle Stage)**

- 41 -

## Decoupled Software Pipelining

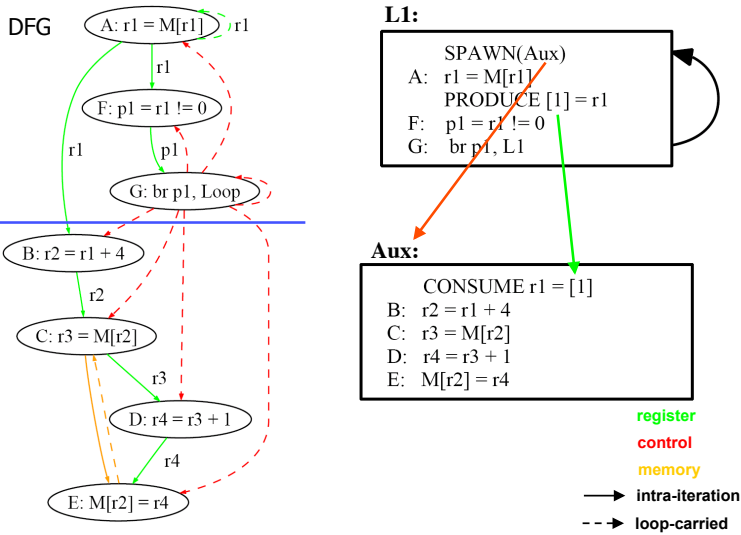# Decoupled Software Pipelining (DSWP)

```
A: while(node)
B:    ncost = doit(node);
C:    cost += ncost;
D:    node = node->next;
```

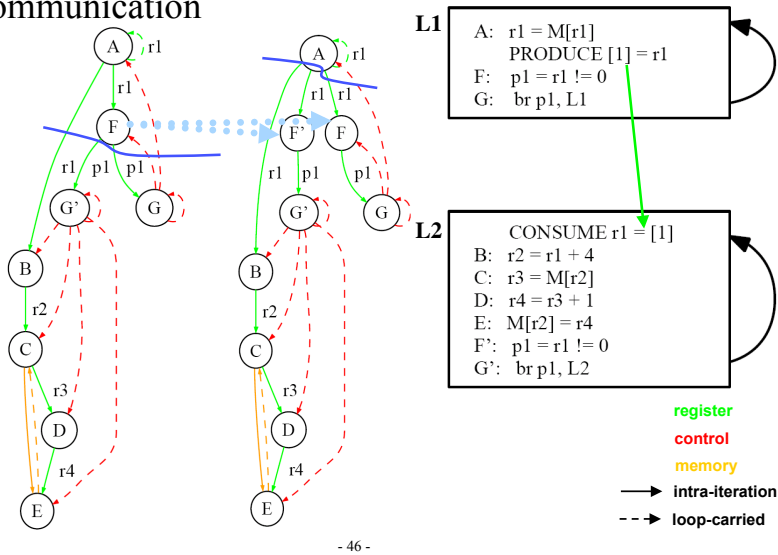**Dependence Graph**



**DAG$_{SCC}$**



**Thread 1**          **Thread 2**



**register**
**control**

→ intra-iteration
-→ loop-carried
■ communication queue

**Inter-thread communication latency is a one-time cost**

- 43 -

# Implementing DSWP

**DFG**



**L1:**
```
        SPAWN(Aux)
A:    r1 = M[r1]
        PRODUCE [1] = r1
F:    p1 = r1 != 0
G:    br p1, L1
```

**Aux:**
```
        CONSUME r1 = [1]
B:    r2 = r1 + 4
C:    r3 = M[r2]
D:    r4 = r3 + 1
E:    M[r2] = r4
```

**register**
**control**
**memory**

→ intra-iteration
-→ loop-carried

- 44 -

# Optimization: Node Splitting
# To Eliminate Cross Thread Control



**L1**
```
A:    r1 = M[r1]
        PRODUCE [1] = r1
F:    p1 = r1 != 0
        PRODUCE [2] = p1
G:    br p1, L1
```

**L2**
```
        CONSUME r1 = [1]
B:    r2 = r1 + 4
C:    r3 = M[r2]
D:    r4 = r3 + 1
E:    M[r2] = r4
        CONSUME p1 = [2]
G':   br p1, L2
```

**register**
**control**
**memory**

→ intra-iteration
-→ loop-carried

- 45 -

# Optimization: Node Splitting To Reduce Communication



**L1**
```
A:   r1 = M[r1]
     PRODUCE [1] = r1
F:   p1 = r1 != 0
G:   br p1, L1
```

**L2**
```
     CONSUME r1 = [1]
B:   r2 = r1 + 4
C:   r3 = M[r2]
D:   r4 = r3 + 1
E:   M[r2] = r4
F':  p1 = r1 != 0
G':  br p1, L2
```

**register**
**control**
**memory**
⟶ **intra-iteration**
⤏ **loop-carried**

# Constraint: Strongly Connected Components

Consider:



Solution: DAG$_{SCC}$

```
         SPAWN(Aux)
A:   r1 = M[r1]
B:   r2 = r1 + 4
C:   r3 = M[r2]
     PRODUCE [1] = r3
     CONSUME r0 = [2]
F:   p1 = r1 != 0
G:   br p1, L1
```

```
     CONSUME r3 = [1]
D:   r4 = r3 + 1
E:   M[r2] = r4
     PRODUCE [2] = r0
```

**Eliminates pipelined/decoupled property**

**register**
**control**
**memory**
⟶ **intra-iteration**
⤏ **loop-carried**

# 2 Extensions to the Basic Transformation

❖ Speculation
  » Break statistically unlikely dependences
  » Form better-balanced pipelines
❖ Parallel Stages
  » Execute multiple copies of certain "large" stages
  » Stages that contain inner loops perfect candidates

# Why Speculation?

```
A: while(node)
B:    ncost = doit(node);
C:    cost += ncost;
D:    node = node->next;
```
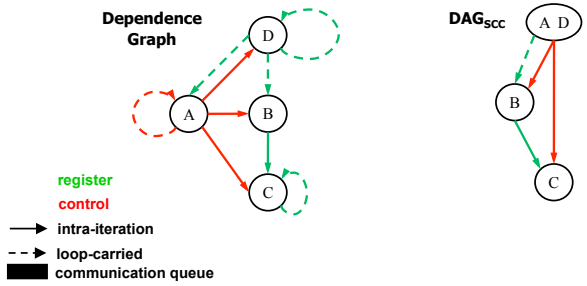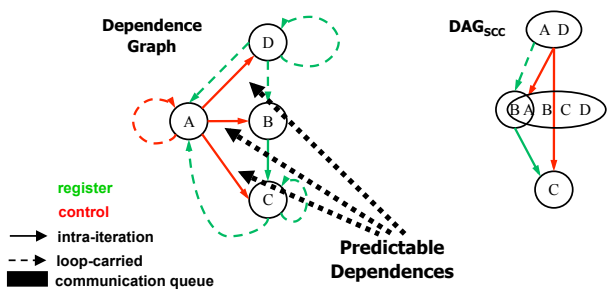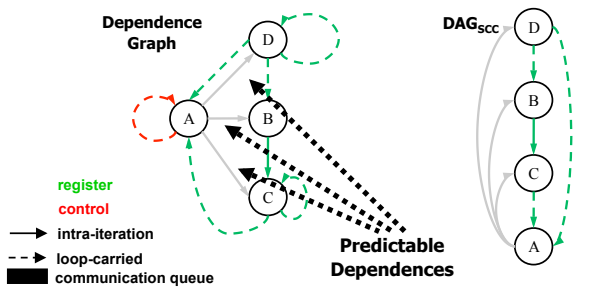
**Dependence Graph**

**DAG$_{SCC}$**

**register**
**control**
→ intra-iteration
- - → loop-carried
■ communication queue

# Why Speculation?

```
A: while(cost < T && node)
B:    ncost = doit(node);
C:    cost += ncost;
D:    node = node->next;
```
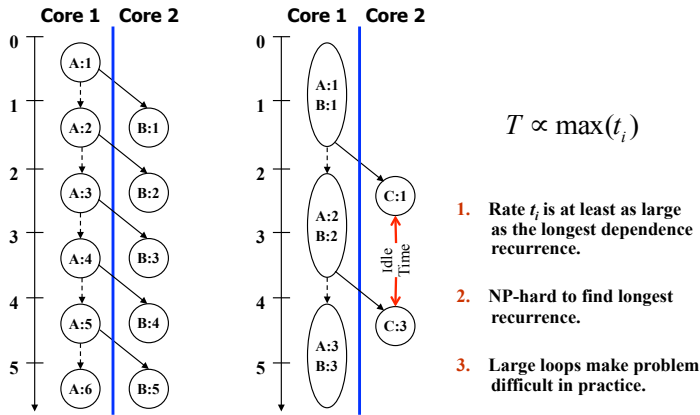
**Dependence Graph**

**DAG$_{SCC}$**

**register**
**control**
→ intra-iteration
- - → loop-carried
■ communication queue

**Predictable Dependences**

# Why Speculation?

```
A: while(cost < T && node)
B:    ncost = doit(node);
C:    cost += ncost;
D:    node = node->next;
```

**Dependence Graph**

**DAG$_{SCC}$**

**register**
**control**
→ intra-iteration
- - → loop-carried
■ communication queue

**Predictable Dependences**

# Execution Paradigm



DAG_SCC with nodes D, B, C, A

register (green)
control (red)
→ intra-iteration
- - → loop-carried
■ communication queue

Misspeculation detected

Misspeculation Recovery
Rerun Iteration 4

# Understanding PMT Performance



$$T \propto \max(t_i)$$

1. Rate $t_i$ is at least as large as the longest dependence recurrence.

2. NP-hard to find longest recurrence.

3. Large loops make problem difficult in practice.

| | |
|---|---|
| Slowest thread: | 1 cycle/iter |
| Iteration Rate: | 1 iter/cycle |

2 cycle/iter
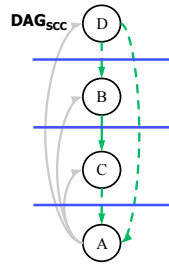
0.5 iter/cycle

# Selecting Dependences To Speculate

```
A: while(cost < T && node)
B:    ncost = doit(node);
C:    cost += ncost;
D:    node = node->next;
```
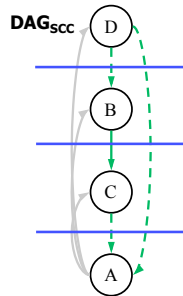


Dependence Graph

register (green)
control (red)
→ intra-iteration
- - → loop-carried
■ communication queue

DAG_SCC with nodes D, B, C, A
Thread 1, Thread 2, Thread 3, Thread 4

# Detecting Misspeculation

```
A¹: while(consume(4))
D :    node = node->next
       produce({0,1},node);
```

**Thread 2**

```
A²: while(consume(5))
B :    ncost = doit(node);
       produce(2,ncost);
D²:    node = consume(0);
```

**Thread 3**

```
A³: while(consume(6))
B³:    ncost = consume(2);
C :    cost += ncost;
       produce(3,cost);
```

**Thread 4**

```
A : while(cost < T && node)
B⁴:    cost = consume(3);
C⁴:    node = consume(1);
       produce({4,5,6},cost < T
                      && node);
```
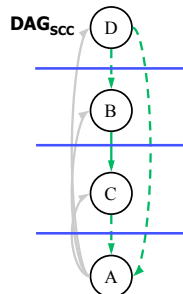
DAG_SCC



- 55 -

# Detecting Misspeculation

**Thread 1**

```
A¹: while(TRUE)
D :    node = node->next
       produce({0,1},node);
```

**Thread 2**

```
A²: while(TRUE)
B :    ncost = doit(node);
       produce(2,ncost);
D²:    node = consume(0);
```

**Thread 3**

```
A³: while(TRUE)
B³:    ncost = consume(2);
C :    cost += ncost;
       produce(3,cost);
```

**Thread 4**

```
A : while(cost < T && node)
B⁴:    cost = consume(3);
C⁴:    node = consume(1);
       produce({4,5,6},cost < T
                      && node);
```
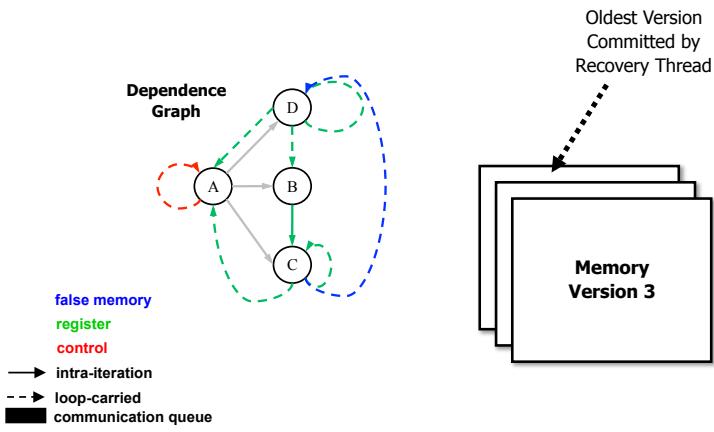
DAG_SCC



- 56 -

# Detecting Misspeculation

**Thread 1**

```
A¹: while(TRUE)
D :    node = node->next
       produce({0,1},node);
```

**Thread 2**

```
A²: while(TRUE)
B :    ncost = doit(node);
       produce(2,ncost);
D²:    node = consume(0);
```

**Thread 3**

```
A³: while(TRUE)
B³:    ncost = consume(2);
C :    cost += ncost;
       produce(3,cost);
```

**Thread 4**

```
A : while(cost < T && node)
B⁴:    cost = consume(3);
C⁴:    node = consume(1);
       if(!(cost < T && node))
          FLAG_MISSPEC();
```
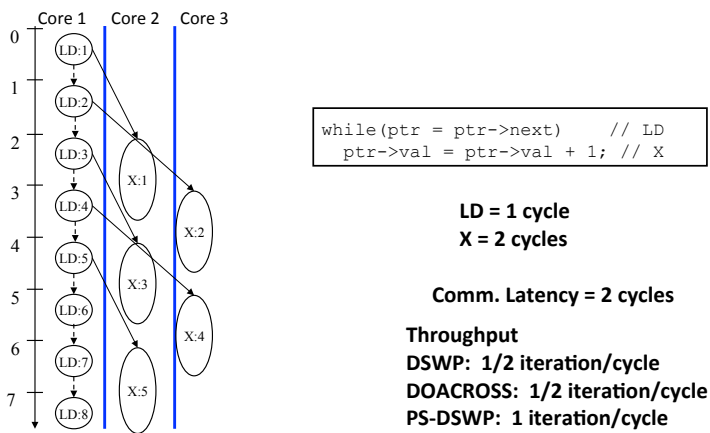
DAG_SCC



- 57 -

# Breaking False Memory Dependences



Dependence Graph

D, A, B, C

Oldest Version Committed by Recovery Thread

Memory Version 3

**false memory**
**register**
**control**

→ **intra-iteration**
- - → **loop-carried**
■ **communication queue**

# Adding Parallel Stages to DSWP



Core 1    Core 2    Core 3

```
while(ptr = ptr->next)      // LD
  ptr->val = ptr->val + 1; // X
```

**LD = 1 cycle**
**X = 2 cycles**

**Comm. Latency = 2 cycles**

**Throughput**
**DSWP:  1/2 iteration/cycle**
**DOACROSS:  1/2 iteration/cycle**
**PS-DSWP:  1 iteration/cycle**

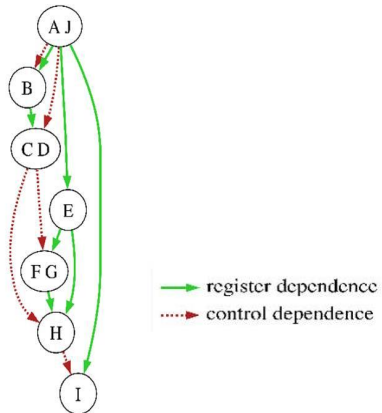# Thread Partitioning

```
        p = list;
      sum = 0;
A:   while (p != NULL) {
B:     id = p->id;
E:       q = p->inner_list;
C:     if (!visited[id]) {
D:       visited[id] = true;
F:       while (foo(q))
G:         q = q->next;
H:       if (q != NULL)
I:         sum += p->value;
       }
J:    p = p->next;
      }
```
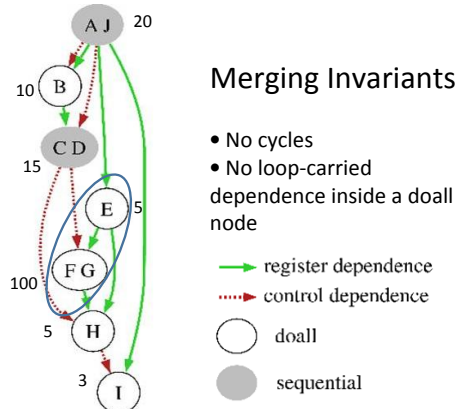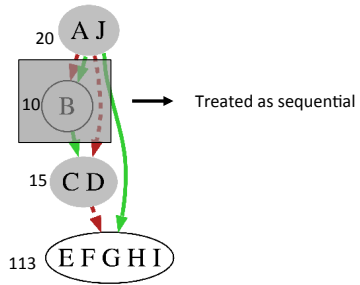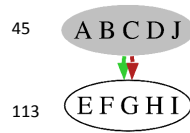


→ register dependence
····→ control dependence

Reduction

# Thread Partitioning: DAG$_{SCC}$

A J
B
C D
E
F G
H
I

register dependence
control dependence

# Thread Partitioning

A J 20
10 B
C D
15
E 5
F G
100
5 H
3 I

## Merging Invariants

• No cycles
• No loop-carried dependence inside a doall node

register dependence
control dependence
doall
sequential

# Thread Partitioning

20 A J
10 B → Treated as sequential
15 C D
113 E F G H I

## Thread Partitioning



- ❖ Modified MTCG[Ottoni, MICRO'05] to generate code from partition

## Discussion Point 1 – Speculation

- ❖ How do you decide what dependences to speculate?
  - » Look solely at profile data?
  - » How do you ensure enough profile coverage?
  - » What about code structure?
  - » What if you are wrong?  Undo speculation decisions at run-time?
- ❖ How do you manage speculation in a pipeline?
  - » Traditional definition of a transaction is broken
  - » Transaction execution spread out across multiple cores

## Discussion Point 2 – Pipeline Structure

- ❖ When is a pipeline a good/bad choice for parallelization?

- ❖ Is pipelining good or bad for cache performance?
  - » Is DOALL better/worse for cache?

- ❖ Can a pipeline be adjusted when the number of available cores increases/decreases?

# CFGs, PCs, and Cross-Iteration Deps

```
1.  r1 = 10
```

```
1.  r1 = r1 + 1

2.  r2 = MEM[r1]

3.  r2 = r2 + 1

4.  MEM[r1] = r2

5.  Branch r1 < 1000
```

No register live outs

# Loop-Level Parallelization: DOALL

```
1.  r1 = 10
```
```
1.  r1 = 9
```
```
1.  r1 = 10
```

```
1.  r1 = r1 + 1

2.  r2 = MEM[r1]

3.  r2 = r2 + 1

4.  MEM[r1] = r2

5.  Branch r1 < 1000
```
```
1.  r1 = r1 + 2

2.  r2 = MEM[r1]

3.  r2 = r2 + 1

4.  MEM[r1] = r2

5.  Branch r1 < 999
```
```
1.  r1 = r1 + 2

2.  r2 = MEM[r1]

3.  r2 = r2 + 1

4.  MEM[r1] = r2

5.  Branch r1 < 1000
```

No register live outs

# Another Example

```
1.  r1 = 10
```

```
1.  r1 = r1 + 1

2.  r2 = MEM[r1]

3.  r2 = r2 + 1

4.  MEM[r1] = r2

5.  Branch r2 == 10
```

No register live outs

## Another Example

| 1. r1 = 10 | 1. r1 = 9 | 1. r1 = 10 |
|---|---|---|
| 1. r1 = r1 + 1 | 1. r1 = r1 + 2 | 1. r1 = r1 + 2 |
| 2. r2 = MEM[r1] | 2. r2 = MEM[r1] | 2. r2 = MEM[r1] |
| 3. r2 = r2 + 1 | 3. r2 = r2 + 1 | 3. r2 = r2 + 1 |
| 4. MEM[r1] = r2 | 4. MEM[r1] = r2 | 4. MEM[r1] = r2 |
| 5. Branch r2 == 10 | 5. Branch r2 == 10 | 5. Branch r2 == 10 |

No register live outs

70

## Speculation

| 1. r1 = 9 | 1. r1 = 10 |
|---|---|
| 1. r1 = r1 + 2 | 1. r1 = r1 + 2 |
| 2. r2 = MEM[r1] | 2. r2 = MEM[r1] |
| 3. r2 = r2 + 1 | 3. r2 = r2 + 1 |
| 4. MEM[r1] = r2 | 4. MEM[r1] = r2 |
| 5. Branch r2 == 10 | 5. Branch r2 == 10 |

No register live outs

71

## Speculation, Commit, and Recovery

| 1. r1 = 9 | 1. r2 = Receive{1} | 1. r1 = 10 |
|---|---|---|
| 1. r1 = r1 + 2 | 2. Branch r2 != 10 | 1. r1 = r1 + 2 |
| 2. r2 = MEM[r1] | 3. MEM[r1] = r2 | 2. r2 = MEM[r1] |
| 3. r2 = r2 + 1 | 4. r2 = Receive{2} | 3. r2 = r2 + 1 |
| 4. Send{1} r2 | 5. Branch r2 != 10 | 4. MEM[r1] = r2 |
| 5. Jump | 6. MEM[r1] = r2 | 5. Jump |
| | 7. Jump | |

No register live outs

| 1. Kill and Continue |
|---|

72

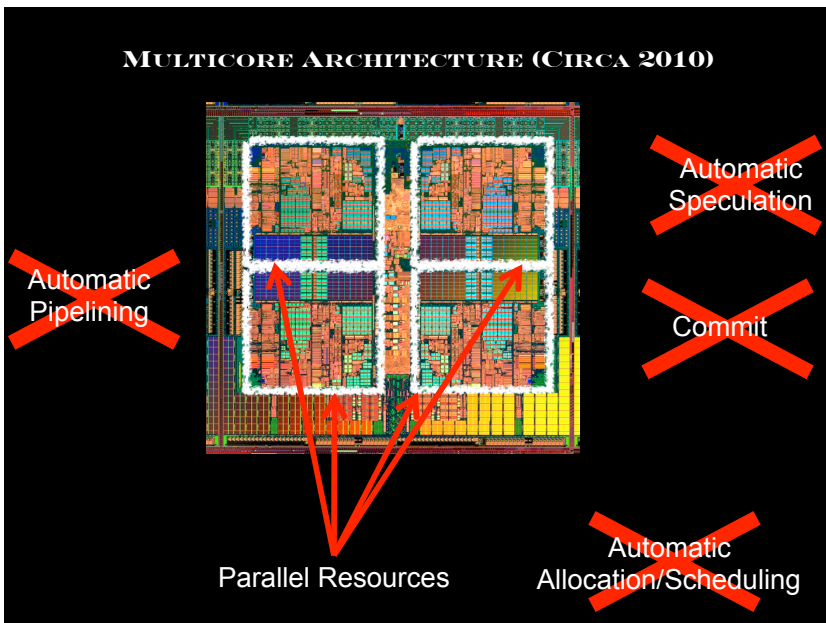## Difficult Dependences

1. r1 = Head

---

1. r1 = MEM[r1]

2. Branch r1 == 0

3. r2 = MEM[r1 + 4]

4. r3 = Work (r2)

5. Print ( r3 )

6. Jump

No register live outs

## DOACROSS

1. r1 = Head

---

1. r1 = MEM[r1]

2. Branch r1 == 0

3. r2 = MEM[r1 + 4]

4. r3 = Work (r2)

5. Print ( r3 )

6. Jump

No register live outs

## PS-DSWP

1. r1 = Head
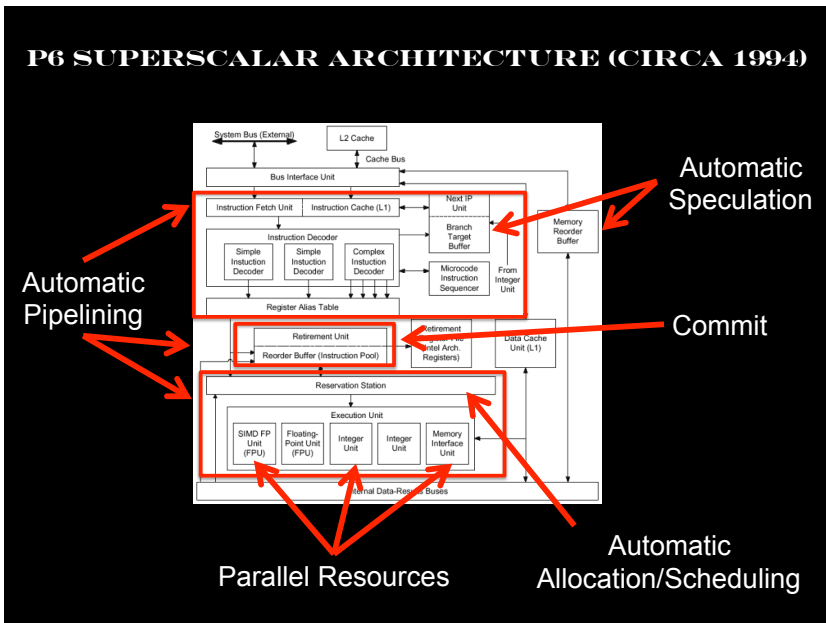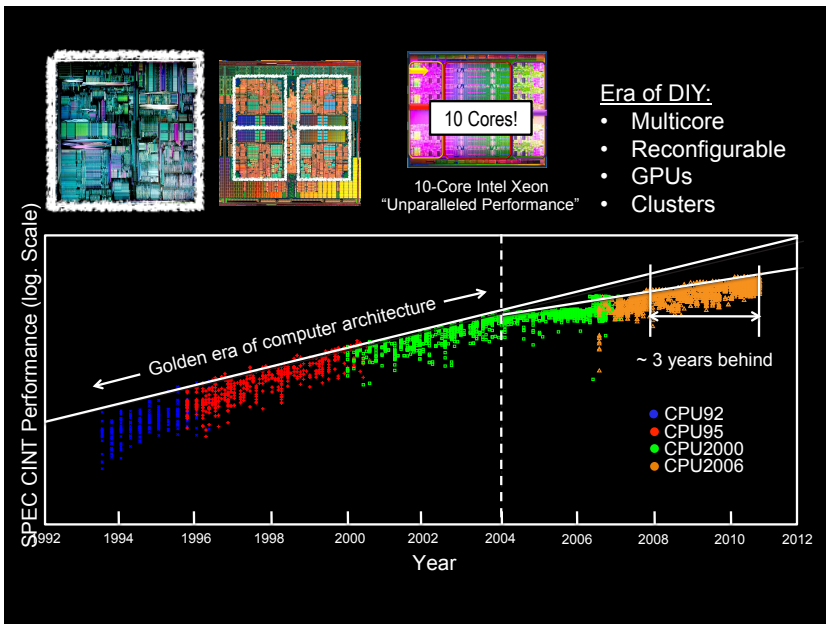
---

1. r1 = MEM[r1]

2. Branch r1 == 0
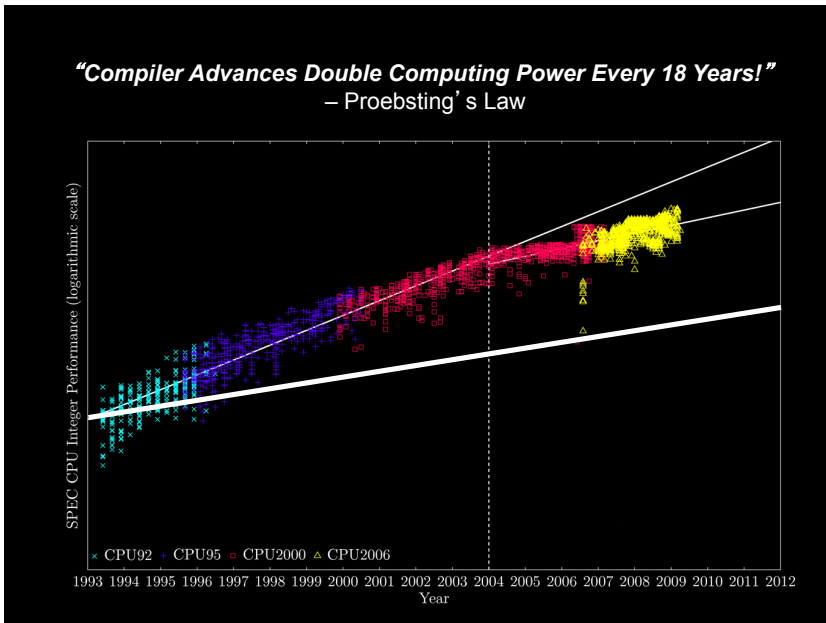
3. r2 = MEM[r1 + 4]

4. r3 = Work (r2)

5. Print ( r3 )

6. Jump

No register live outs

Era of DIY:
- Multicore
- Reconfigurable
- GPUs
- Clusters

10 Cores!

10-Core Intel Xeon
"Unparalleled Performance"

# P6 SUPERSCALAR ARCHITECTURE (CIRCA 1994)



Automatic Speculation

Automatic Pipelining

Commit

Parallel Resources

Automatic Allocation/Scheduling

# MULTICORE ARCHITECTURE (CIRCA 2010)



Automatic Speculation

Automatic Pipelining

Commit

Parallel Resources

Automatic Allocation/Scheduling
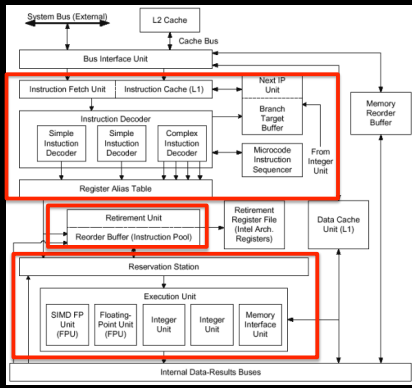
| | | | | | |
|---|---|---|---|---|---|
| ABCPL | CORRELATE | GLU | Mentat | Parafrase2 | pC++ |
| ACE | CPS | GUARD | Legion | Paralation | SCHEDULE |
| ACT++ | CRL | HA-L. | Meta Chaos | Parallel-C++ | SciTL |
| Active messages | CSP | Haskell | Midway | Parallaxis | POET |
| Adl | Cthreads | HPC++ | Millipede | ParC | SDDA. |
| Adsmith | CUMULVS | JAVAR. | CparPar | ParLib++ | SHMEM |
| ADDAP | DAGGER | HORUS | Mirage | ParLin | SIMPLE |
| AFAPI | DAPPLE | HPC | MpC | Parmacs | Sina |
| ALWAN | Data Parallel C | IMPACT | MOSIX | Parti | SISAL. |
| AM | DC++ | ISIS. | Modula-P | pC | distributed smalltalk |
| AMDC | DCE++ | JAVAR | Modula-2* | pC++ | SML |
| AppLeS | DDD | JADE | Multipol | PCN | SONiC |
| Amoeba | DICE. | Java RMI | MPI | PCP: | Split-C. |
| ARTS | DIPC | javaPG | MPC++ | PH | SR |
| Athapascan-0b | DOLIB | JavaSpace | Munin | PEACE | Sthreads |
| Aurora | DOME | JIDL | Nano-Threads | PCU | Strand. |
| Automap | DOSMOS. | Joyce | NESL | PET | SUIF. |
| bb_threads | DRL | Khoros | NetClasses++ | PETSc | Synergy |
| Blaze | DSM-Threads | Karma | Nexus | PENNY | Telegrphos |
| BSP | Ease . | KOAN/Fortran-S | Nimrod | Phosphorus | SuperPascal |
| BlockComm | ECO | LAM | NOW | POET. | TCGMSG. |
| C*. | Eiffel | Lilac | Objective Linda | Polaris | Threads.h++. |
| "C* in C | Eilean | Linda | Occam | POOMA | TreadMarks |
| C** | Emerald | JADA | Omega | POOL-T | TRAPPER |
| CarlOS | EPL | WWWinda | OpenMP | PRESTO | uC++ |
| Cashmere | Excalibur | ISETL-Linda | Orca | P-RIO | UNITY |
| C4 | Express | ParLin | OOF90 | Prospero | UC |
| CC++ | Falcon | Eilean | P++ | Proteus | V |
| Chu | Filaments | P4-Linda | P3L | QPC++ | ViC* |
| Charlotte | FM | Glenda | p4-Linda | PVM | Visifold V-NUS |
| Charm | FLASH | POSYBL | Pablo | PSI | VPE |
| Charm++ | The FORCE | Objective-Linda | PADE | PSDM | Win32 threads |
| Cid | Fork | LiPS | PADRE | Quake | WinPar |
| Cilk | Fortran-M | Locust | Panda | Quark | WWWinda |
| CM-Fortran | FX | Lparx | Papers | Quick Threads | XENOOPS |
| Converse | GA | Lucid | AFAPI | Sage++ | XPC |
| Code | GAMMA | Maisie | Para++ | SCANDAL | Zounds |
| COOL | Glenda | Manifold | Paradigm | SAM | ZPL |



Parallel Library Calls

Threads

Time

Realizable parallelism

Threads

Time

Credit: Jack Dongarra



*"Compiler Advances Double Computing Power Every 18 Years!"*
– Proebsting's Law

SPEC CPU Integer Performance (logarithmic scale)

× CPU92   + CPU95   □ CPU2000   △ CPU2006

1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012
Year

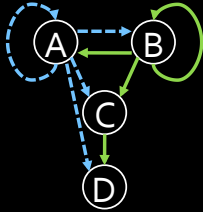# P6 SUPERSCALAR ARCHITECTURE



# Spec-PS-DSWP
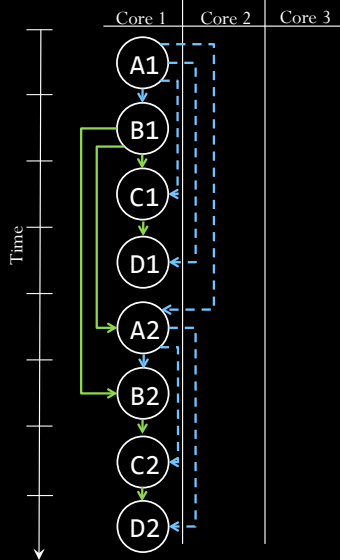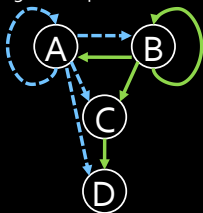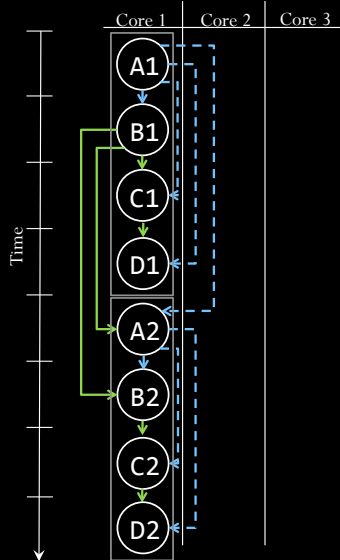


---

Example

```
A: while (node) {
B:    node = node->next;
C:    res = work(node);
D:    write(res);
   }
```

Program Dependence Graph



- - - → Control Dependence
———→ Data Dependence
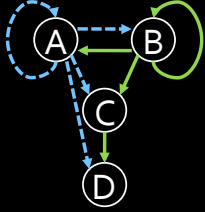
---

# Spec-DOALL

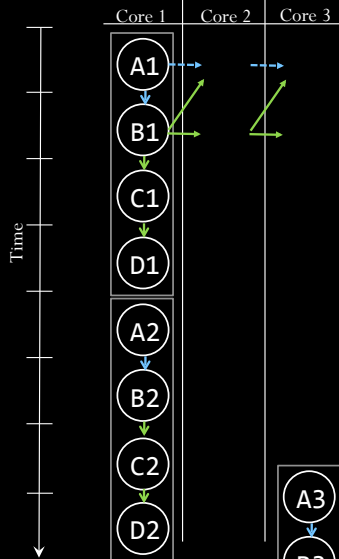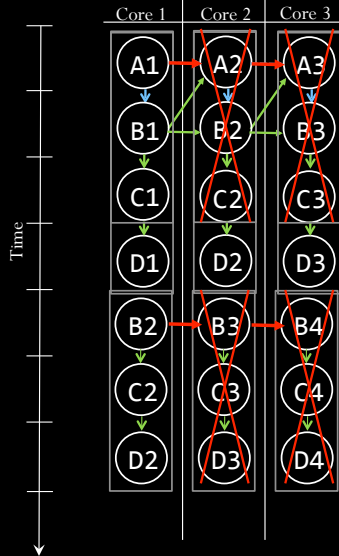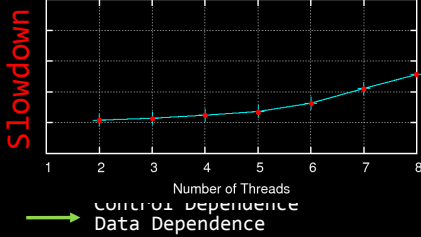Example

```
A: while (node) {
B:    node = node->next;
C:    res = work(node);
D:    write(res);
   }
```

Program Dependence Graph



- - - → Control Dependence
———→ Data Dependence

Example

```
A: while (node) {
B:    node = node->next;
C:    res = work(node);
D:    write(res);
   }
```

Program Dependence Graph

- - - → Control Dependence
──────→ Data Dependence

Core 1 | Core 2 | Core 3
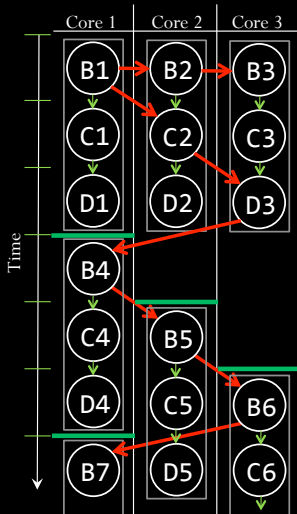
A1 B1 C1 D1 A2 B2 C2 D2 A3

---

Example

```
A: while (node) {
B:    node = node->next;
C:    res = work(node);
D:    write(res);
   }
```

Program Dependence Graph

197.parser

Slowdown

Number of Threads
1  2  3  4  5  6  7  8

Data Dependence

Core 1 | Core 2 | Core 3

A1 A2 A3
B1 B2 B3
C1 C2 C3
D1 D2 D3
B2 B3 B4
C2 C3 C4
D2 D3 D4

---

Spec-DOACROSS
Throughput: 1 iter/cycle

Core 1 | Core 2 | Core 3

B1 B2 B3
C1 C2 C3
D1 D2 D3
B4
C4 B5
D4 C5 B6
B7 D5 C6

Spec-DSWP
Throughput: 1 iter/cycle

Core 1 | Core 2 | Core 3

B1
B2
B3 C2
B4 C3 D2
B5 C4 D3
B6 C5 D4
B7 C6 D5

Time

## Comparison: Spec-DOACROSS and Spec-DSWP

```
Comm.Latency = 1: 1 iter/cycle    Comm.Latency = 1: 1 iter/cycle
Comm.Latency = 2: 0.5 iter/cycle  Comm.Latency = 2: 1 iter/cycle
```



## Spec-DOACROSS vs. Spec-DSWP

[MICRO 2010]

Geomean of 11 benchmarks on the same _cluster_

Performance Speedup (X)

- TLS
- Spec-PS-DSWP

(Number of Total Cores, Number of Nodes)



## Performance relative to Best Sequential

128 Cores in 32 Nodes with Intel Xeon Processors [MICRO 2010]

Performance Speedup (X)

- 052.alvinn
- 130.li
- 164.gzip
- 179.art
- 197.parser
- 256.bzip2
- 456.hmmer
- crc32
- blackscholes
- 464.h264ref
- swaptions

(Number of Total Cores, Number of Nodes)

## CFGs and PCs

94

95

96