# EECS 583 – Automatic Parallelization Via Decoupled Software Pipelining

*University of Michigan*
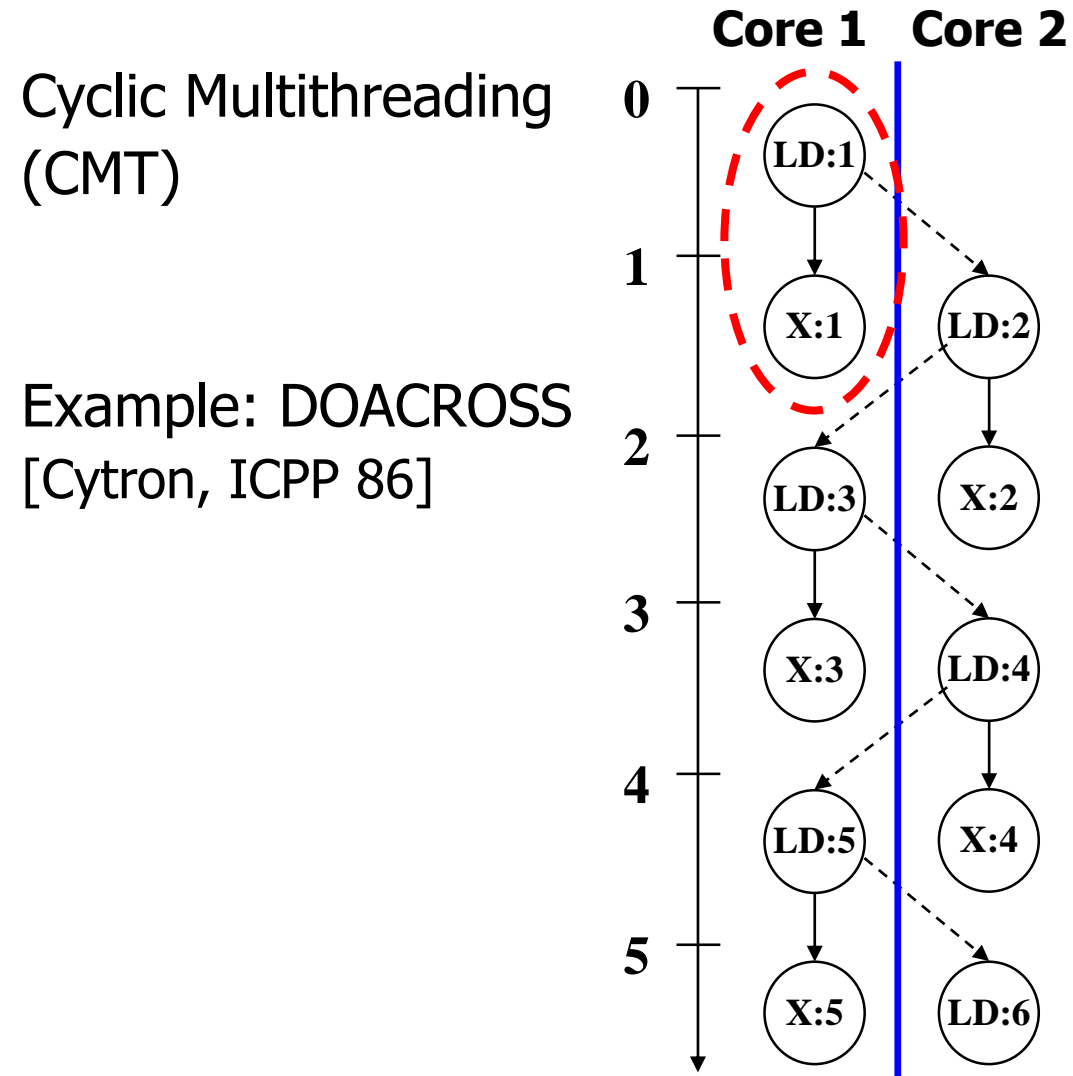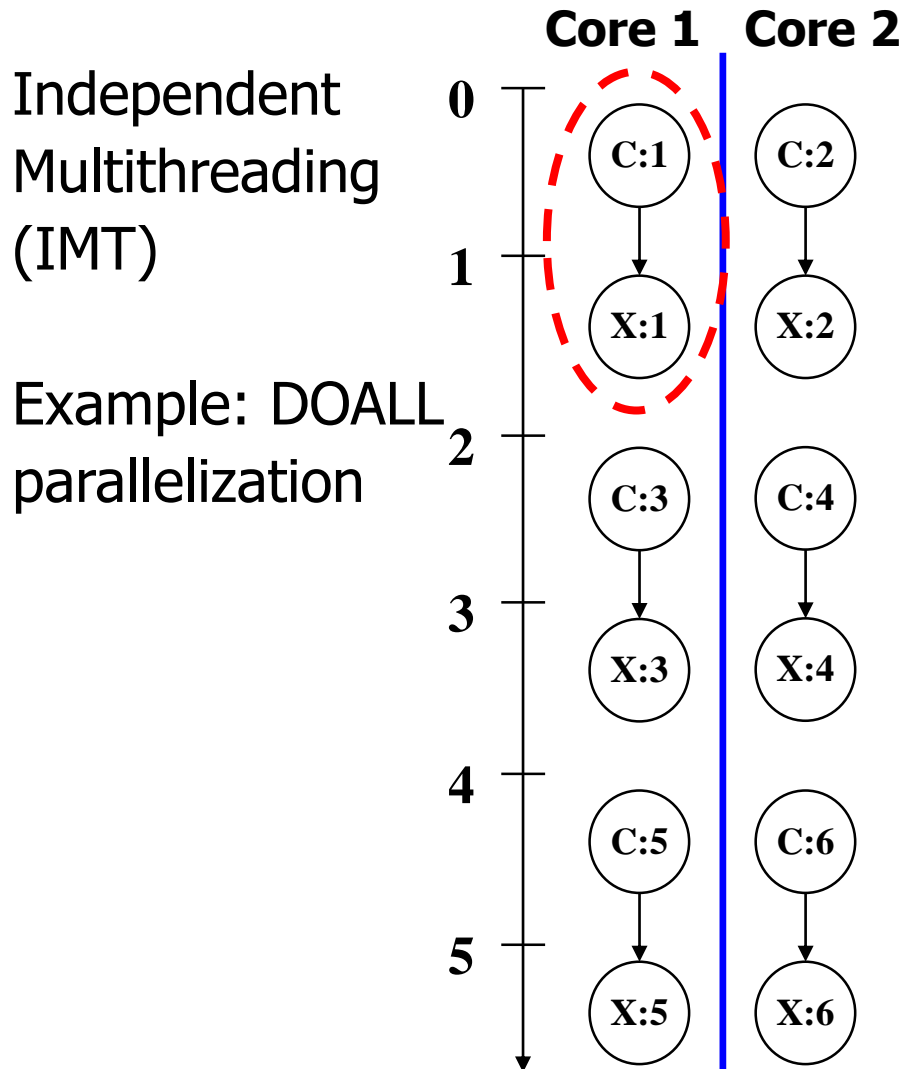
*November 13, 2023*

# Parallelization: Scientific vs Non-Scientific Codes

Scientific Codes (FORTRAN-like)

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

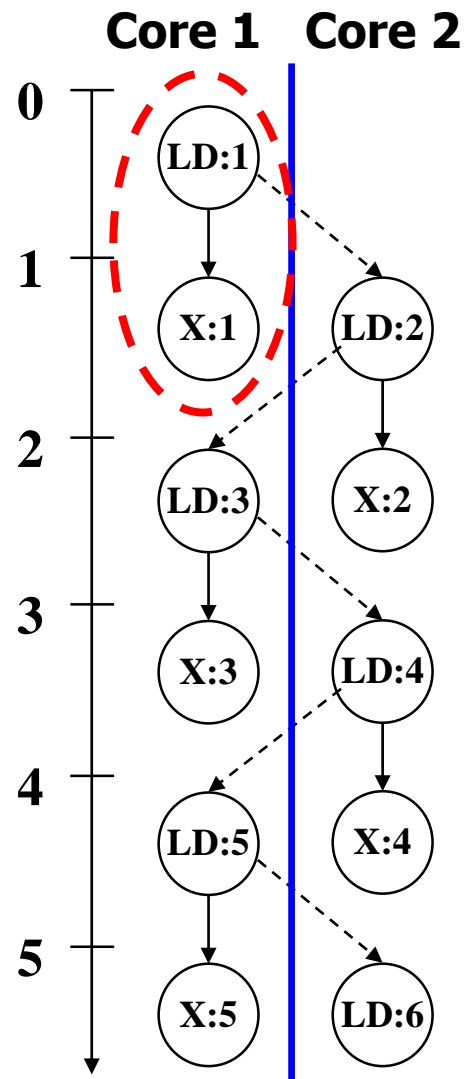General-purpose Codes (legacy C/C++)

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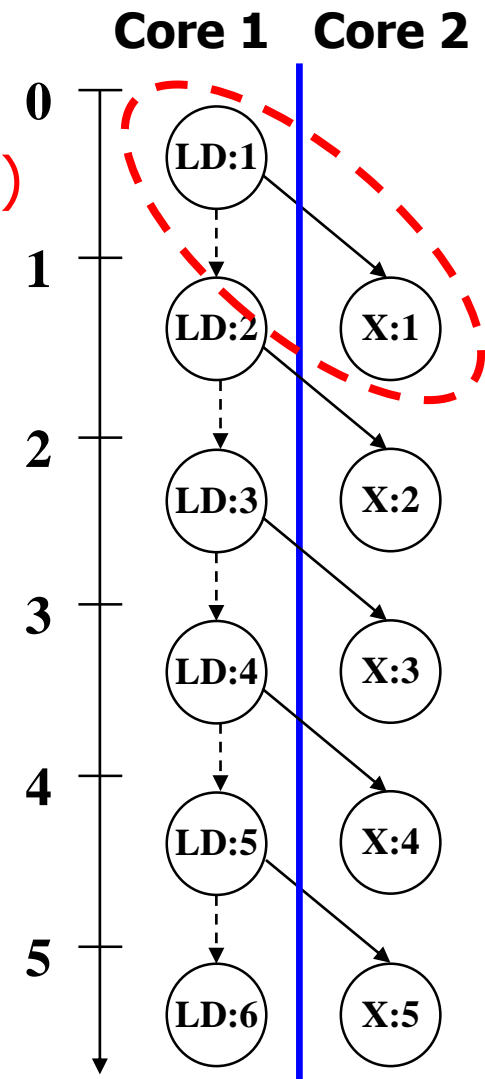
```
while(ptr = ptr->next)     // LD
    ptr->val = ptr->val + 1; // X
```



Cyclic Multithreading (CMT)
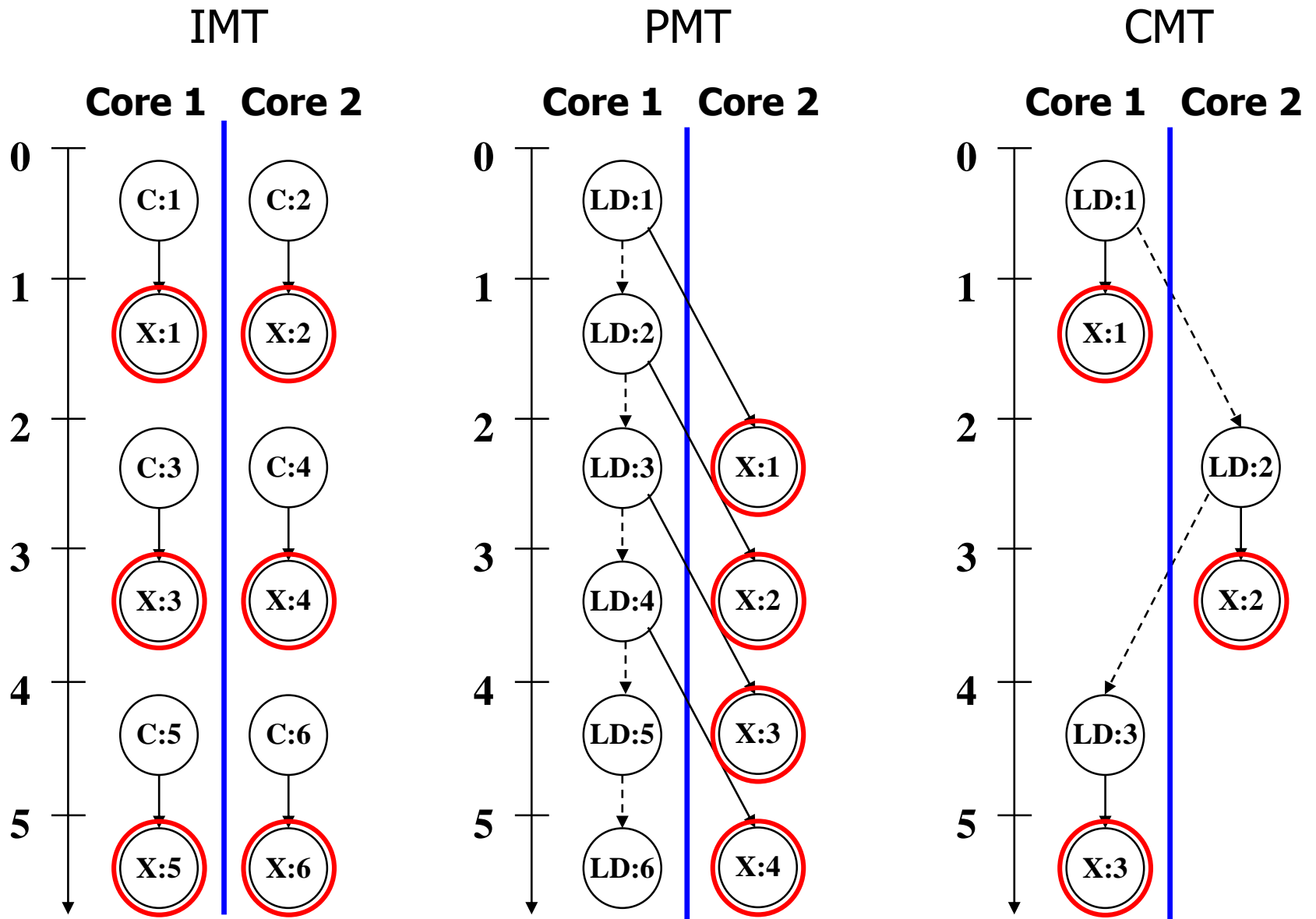
Pipelined Multithreading (PMT)

Example: DSWP
**[PACT 2004]**
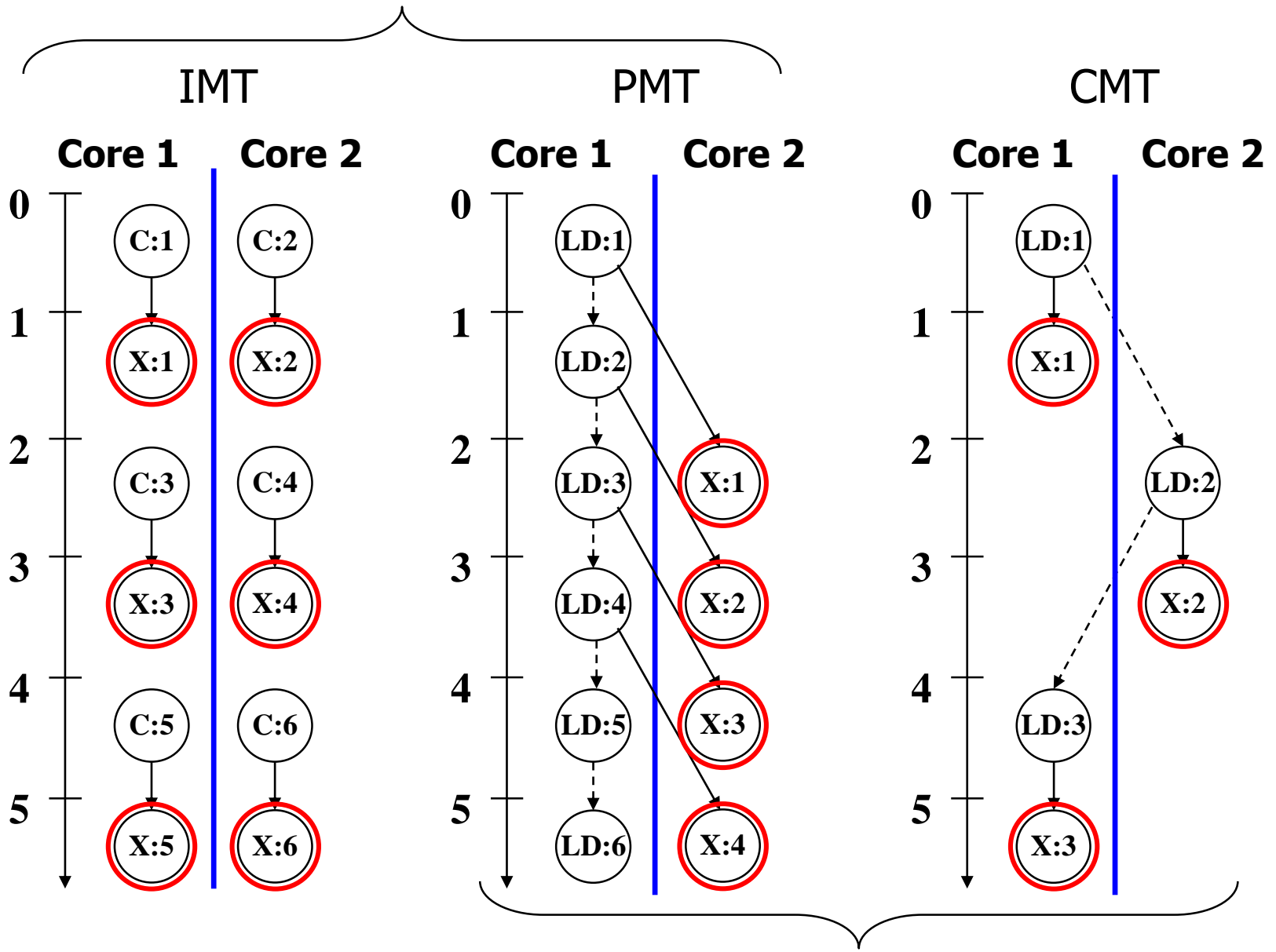
# Comparison: IMT, PMT, CMT



IMT

| | Core 1 | Core 2 |
|---|---|---|
| 0 | C:1 | C:2 |
| 1 | X:1 | X:2 |
| 2 | C:3 | C:4 |
| 3 | X:3 | X:4 |
| 4 | C:5 | C:6 |
| 5 | X:5 | X:6 |

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

PMT

| | Core 1 | Core 2 |
|---|---|---|
| 0 | LD:1 | |
| 1 | LD:2 | |
| 2 | LD:3 | X:1 |
| 3 | LD:4 | X:2 |
| 4 | LD:5 | X:3 |
| 5 | LD:6 | X:4 |

1 iter/cycle
1 iter/cycle

CMT

| | Core 1 | Core 2 |
|---|---|---|
| 0 | LD:1 | |
| 1 | X:1 | |
| 2 | | LD:2 |
| 3 | | X:2 |
| 4 | LD:3 | |
| 5 | X:3 | |

1 iter/cycle
0.5 iter/cycle

# Comparison: IMT, PMT, CMT

**Thread-local Recurrences ➜ Fast Execution**



**Cross-thread Dependences ➜ Wide Applicability**

- 4 -

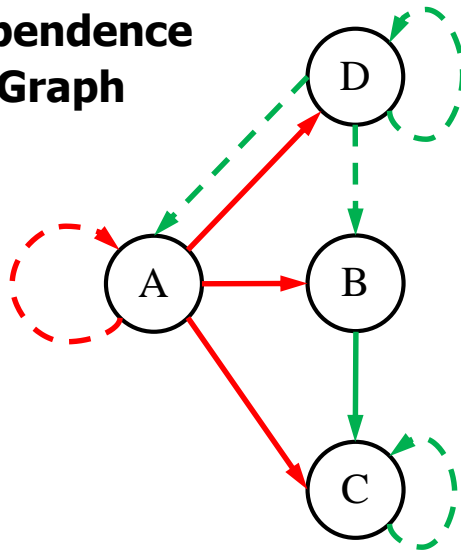# Decoupled Software Pipelining
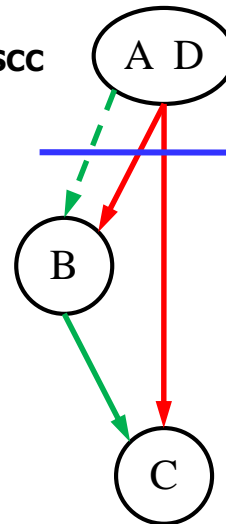
# Decoupled Software Pipelining (DSWP)

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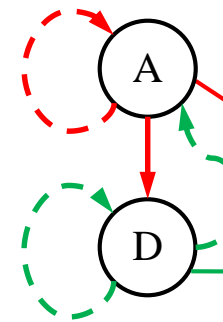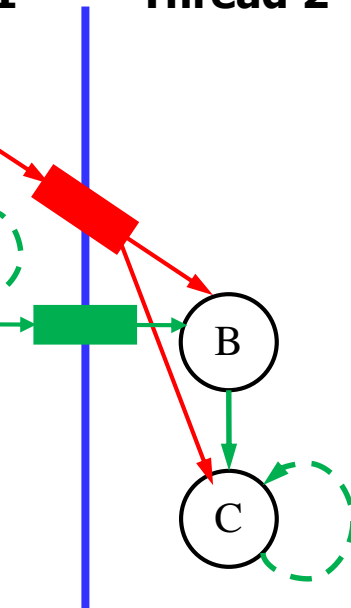
**Dependence Graph**



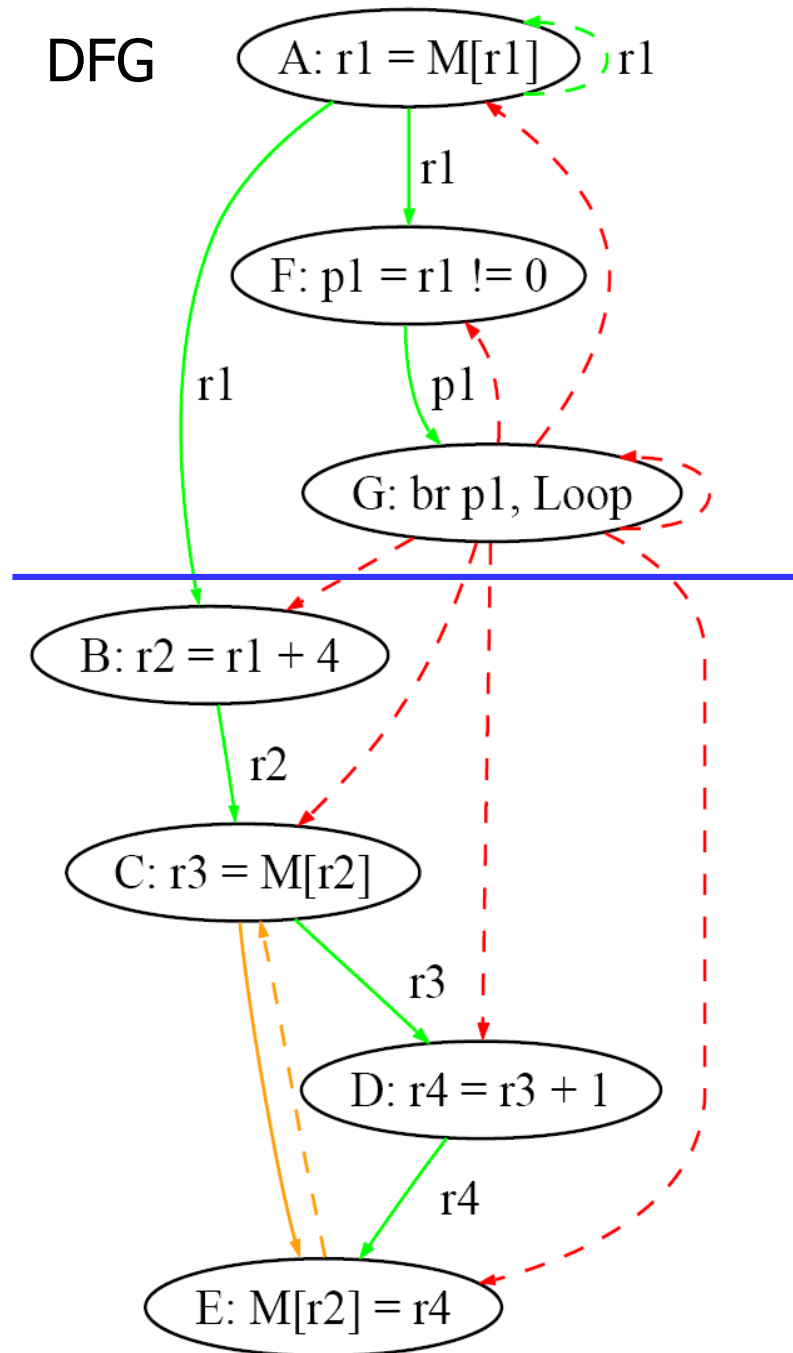**DAG_SCC**

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

register

control

→ intra-iteration

- - -▶ loop-carried

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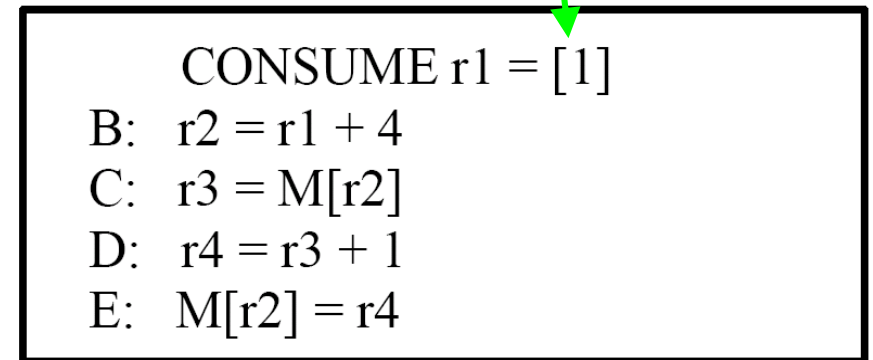
**Inter-thread communication latency is a one-time cost**

# Implementing DSWP



DFG

A: r1 = M[r1]   r1

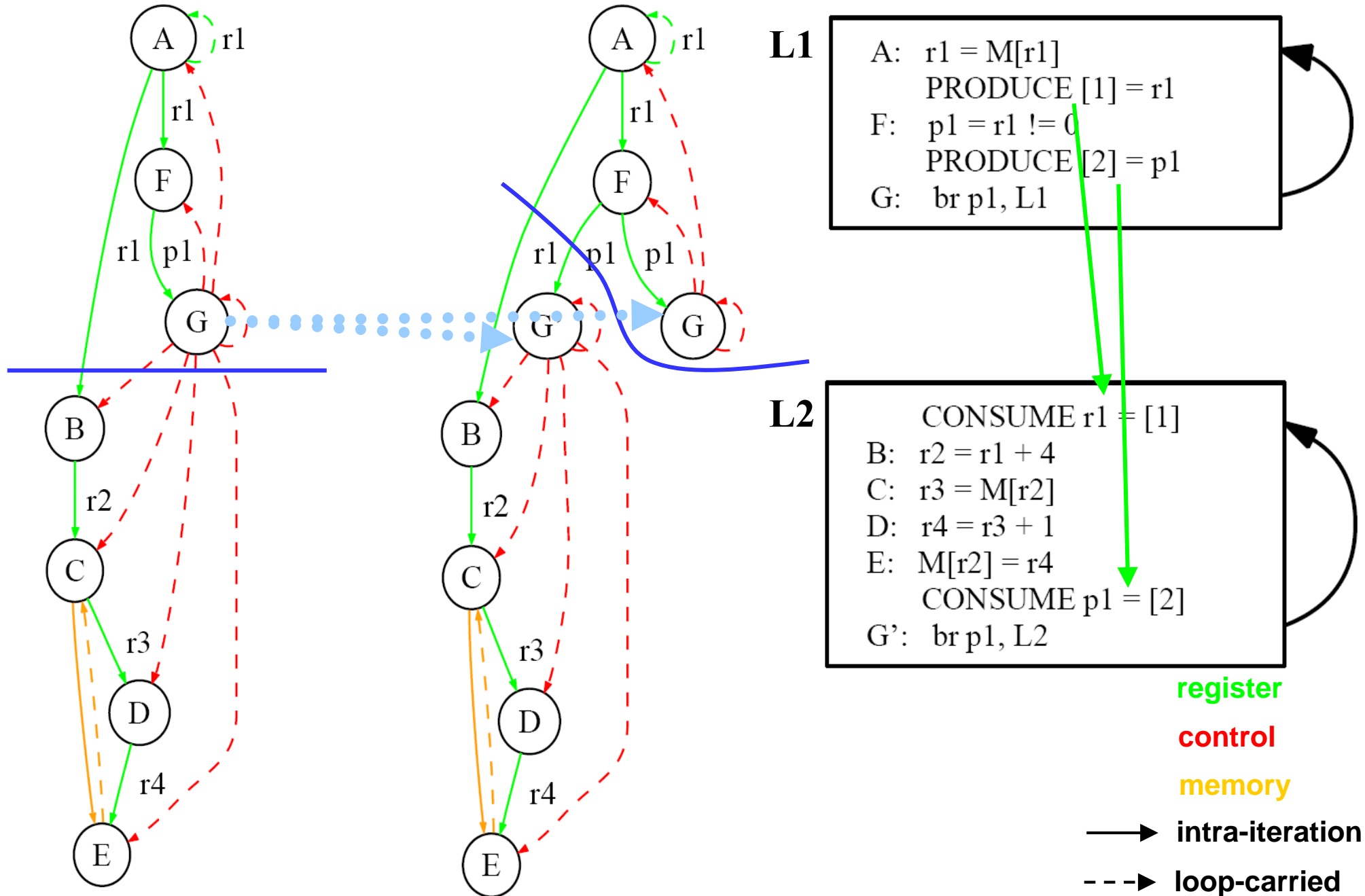r1

F: p1 = r1 != 0

r1      p1

G: br p1, Loop

B: r2 = r1 + 4

r2

C: r3 = M[r2]

r3

D: r4 = r3 + 1

r4

E: M[r2] = r4

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

# 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**

# 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
```
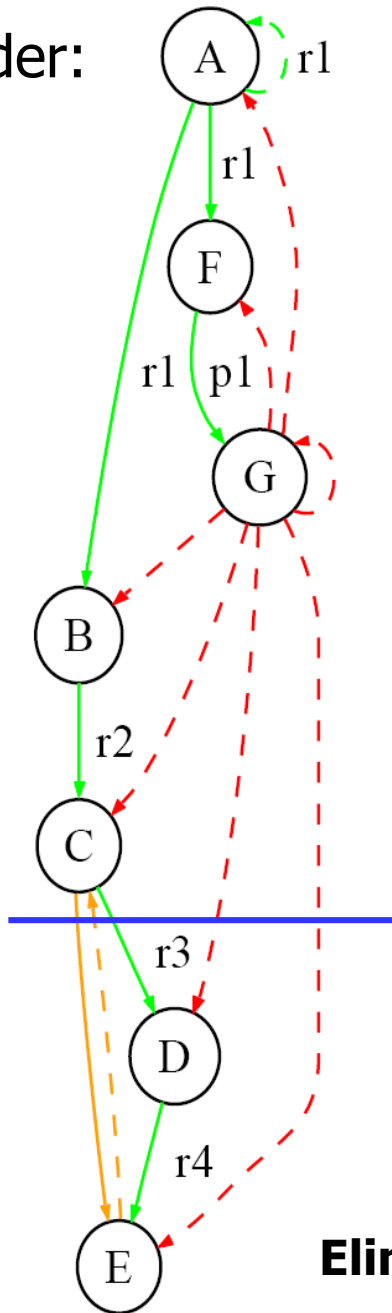
register
control
memory

→ intra-iteration
--→ loop-carried

# Constraint: Strongly Connected Components

Consider:

Solution: DAG$_{SCC}$
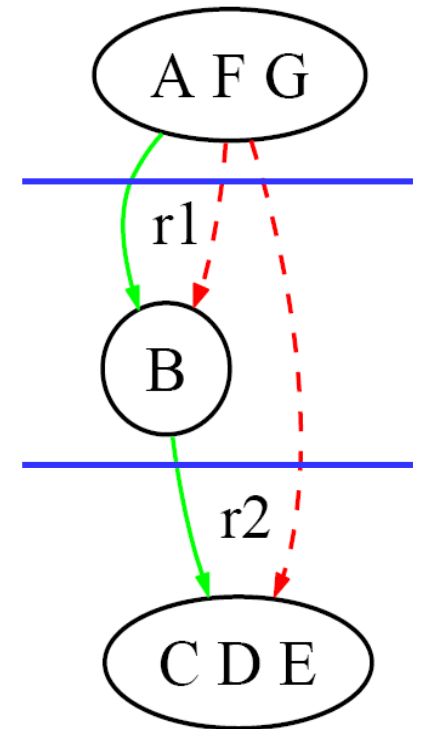


```
        SPAWN(Aux)
A:   r1 = M[r1]
B:   r2 = r1 + 4
C:   r3 = M[r2]
     PRODUCE [1] = r3
     CONSUME r0 = [2]
F:   p1 = r1 != 0
G:   br p1, L1
```

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

**Eliminates pipelined/decoupled property**

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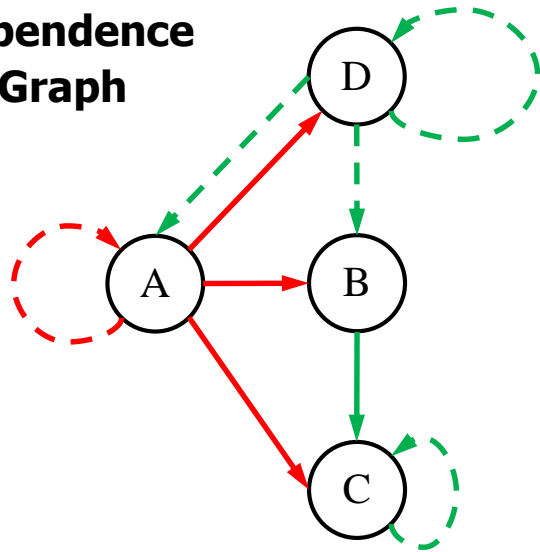**control**
**memory**
⟶ **intra-iteration**
- - -▶ **loop-carried**

# Speculation – Break Statistically Unlikely Dependences

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

**Dependence Graph**
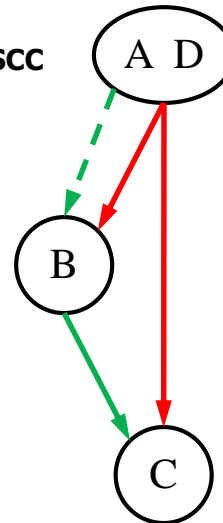
**DAG$_{SCC}$**

**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;
```
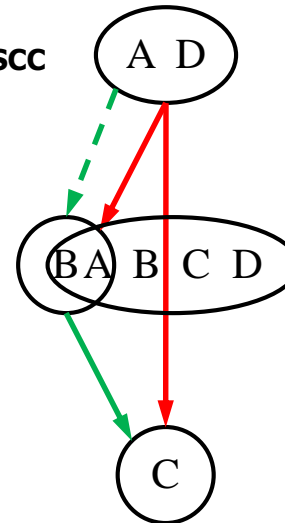
**Dependence Graph**

**DAG<sub>SCC</sub>**

**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**

**register**
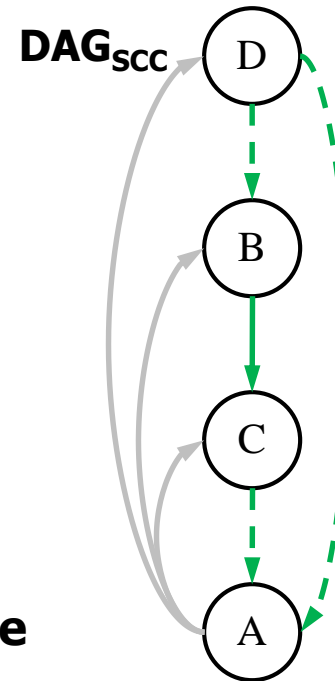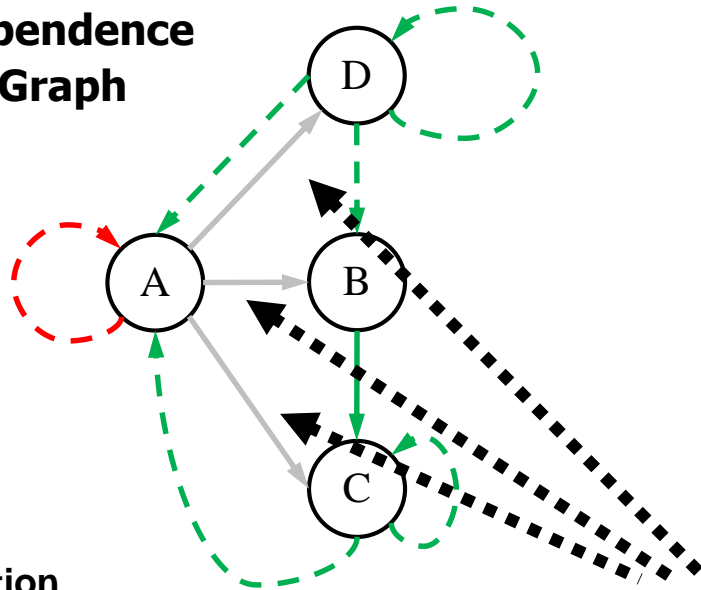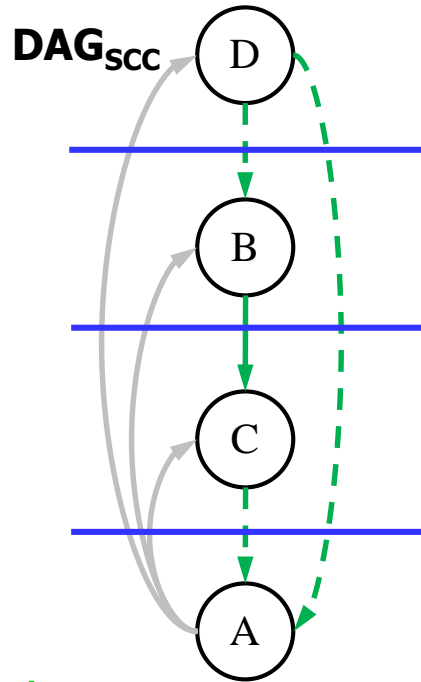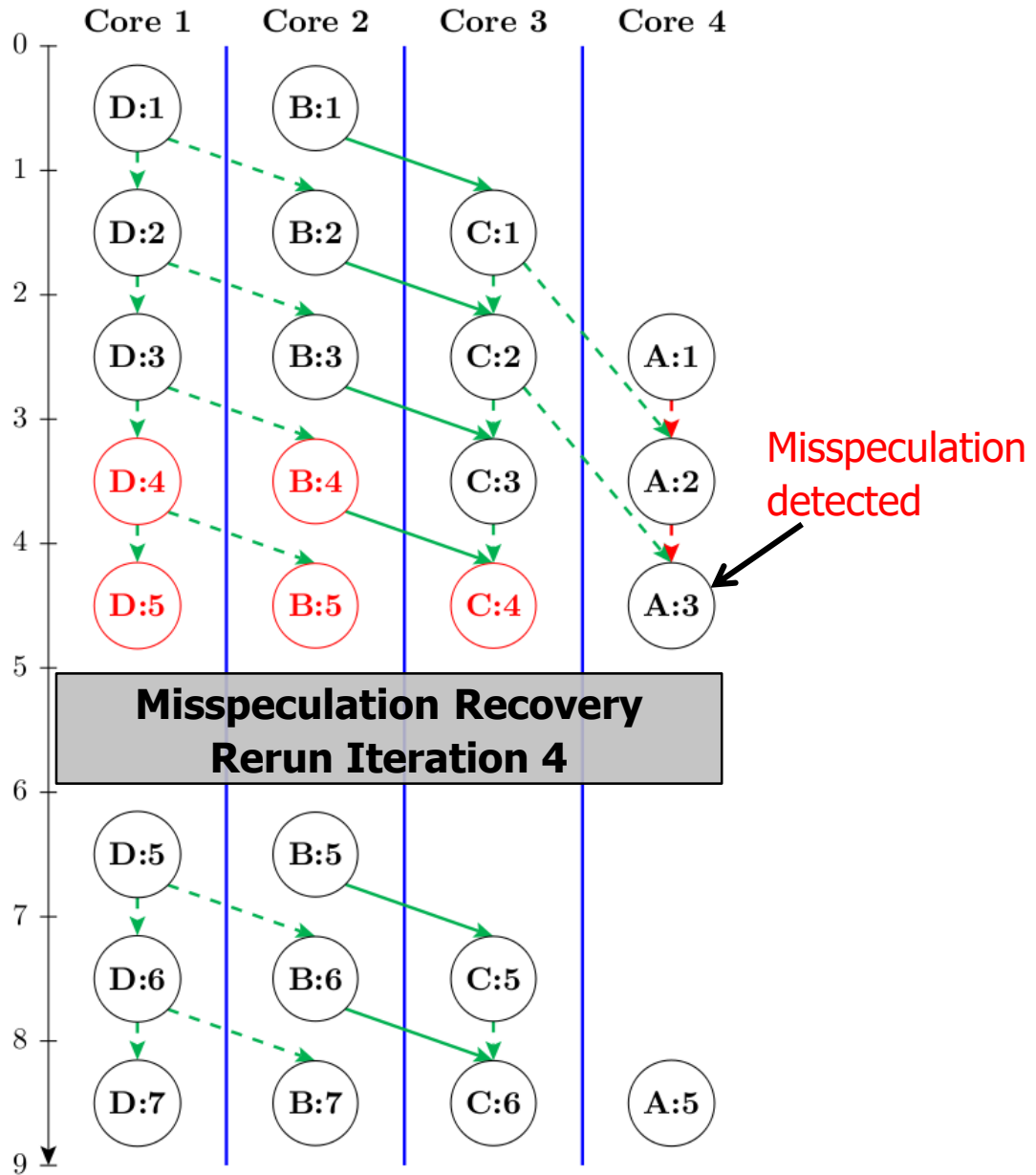
**control**

→ intra-iteration

⇢ loop-carried

▬ communication queue

**Predictable Dependences**

**DAG$_{SCC}$**

# Execution Paradigm

# Evaluation: Dual Core vs Single Core

# References

❖ "Automatic Thread Extraction with Decoupled Software Pipelining," G. Ottoni, R. Rangan, A. Stoler, and D. I. August, *Proceedings of the 38th IEEE/ACM International Symposium on Microarchitecture*, Nov. 2005

❖ "Revisiting the Sequential Programming Model for Multi-Core," M. J. Bridges, N. Vachharajani, Y. Zhang, T. Jablin, and D. I. August, Proc 40th IEEE/ACM International Symposium on Microarchitecture, December 2007.