
Topic 14: Parallelism

COS 320

Compiling Techniques

Princeton University
Spring 2018

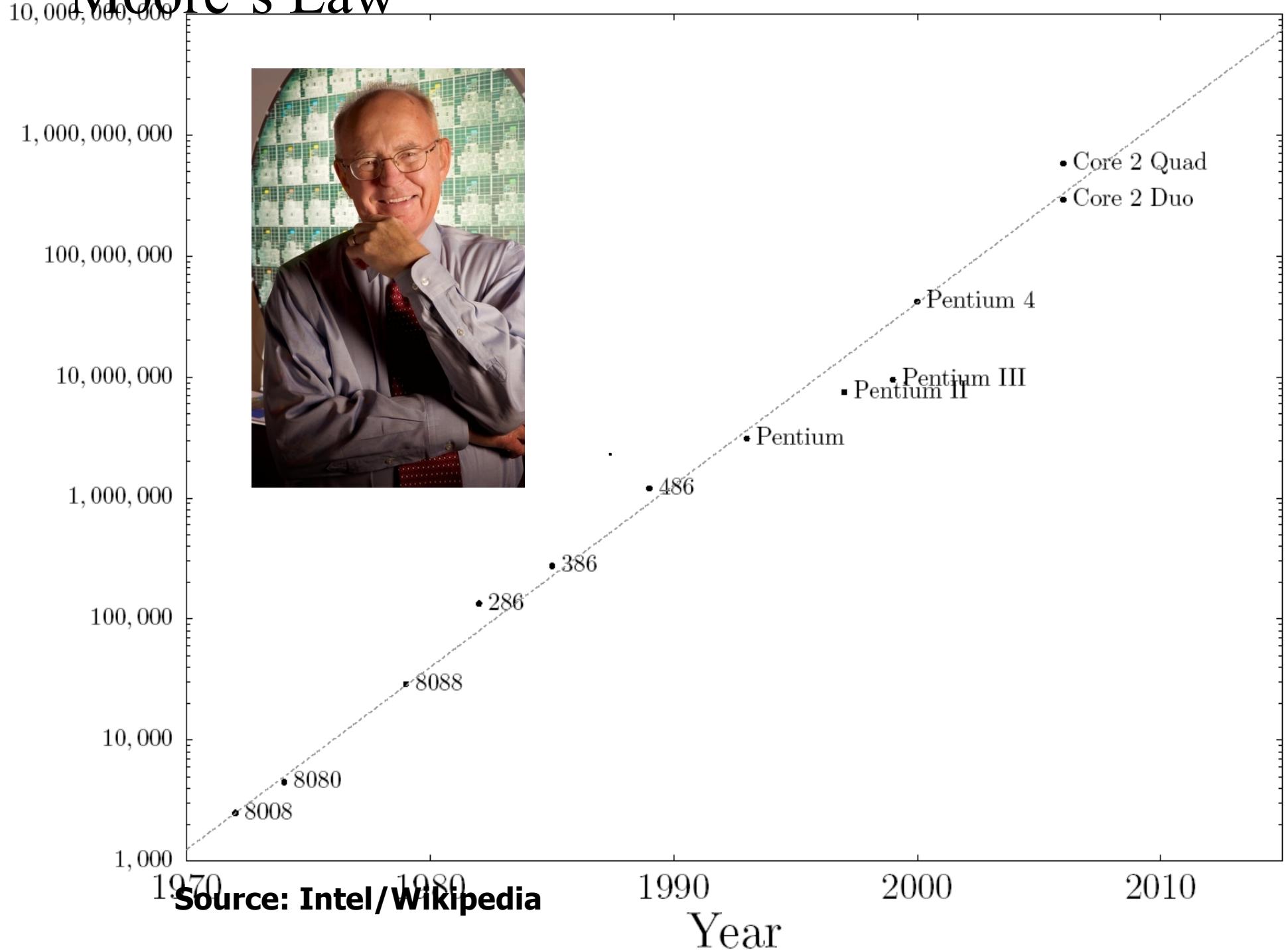
Prof. David August

Final Exam!

- Thursday May 3 in class
- Closed book, closed notes

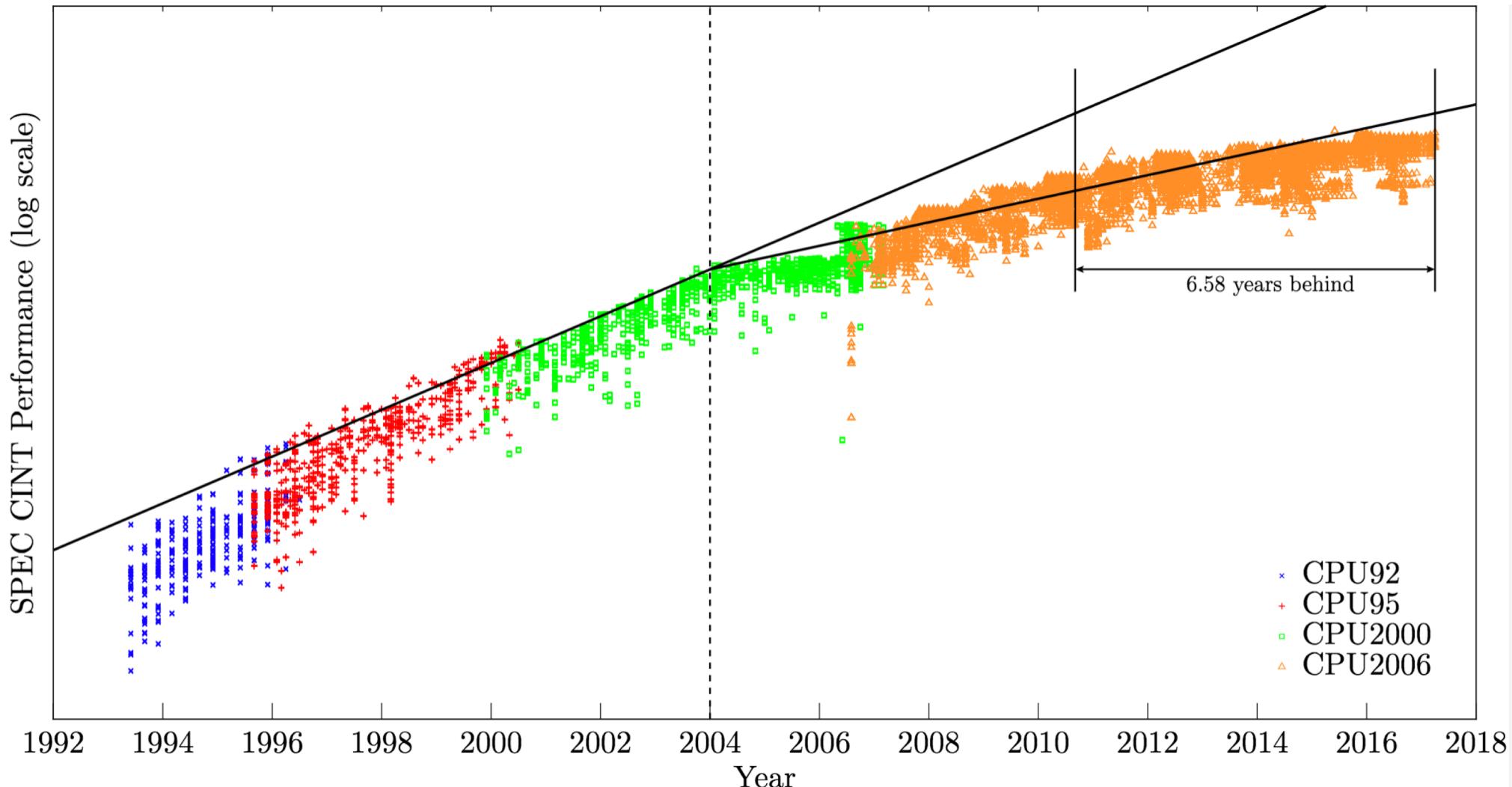
Moore's Law

Transistor Count



Source: Intel/Wikipedia

Single-Threaded Performance Not Improving



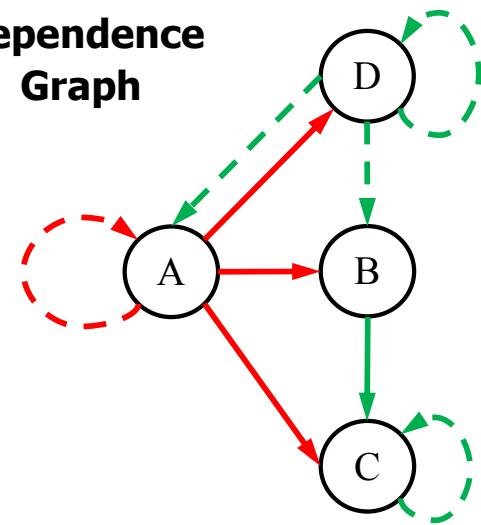
Decoupled Software Pipelining

Decoupled Software Pipelining (DSWP)

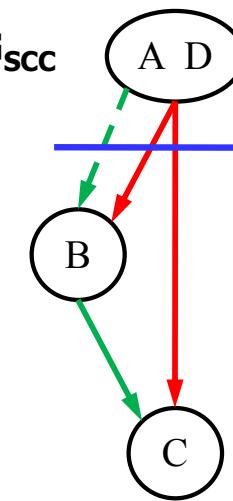
```

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

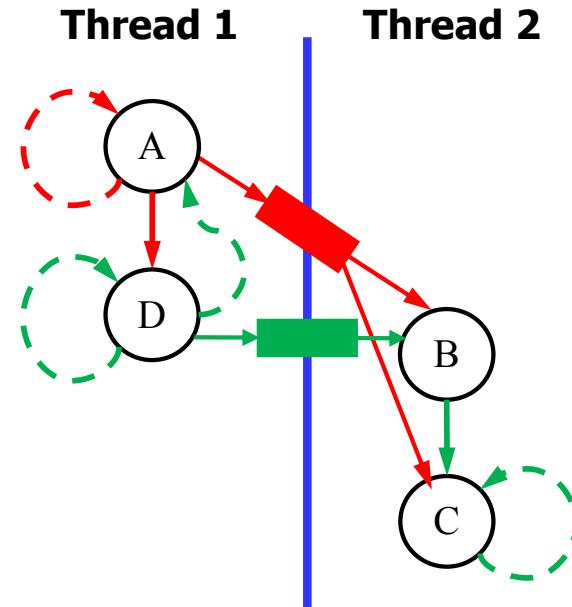
Dependence Graph



DAG_{SCC}



Thread 1



Thread 2

register

control

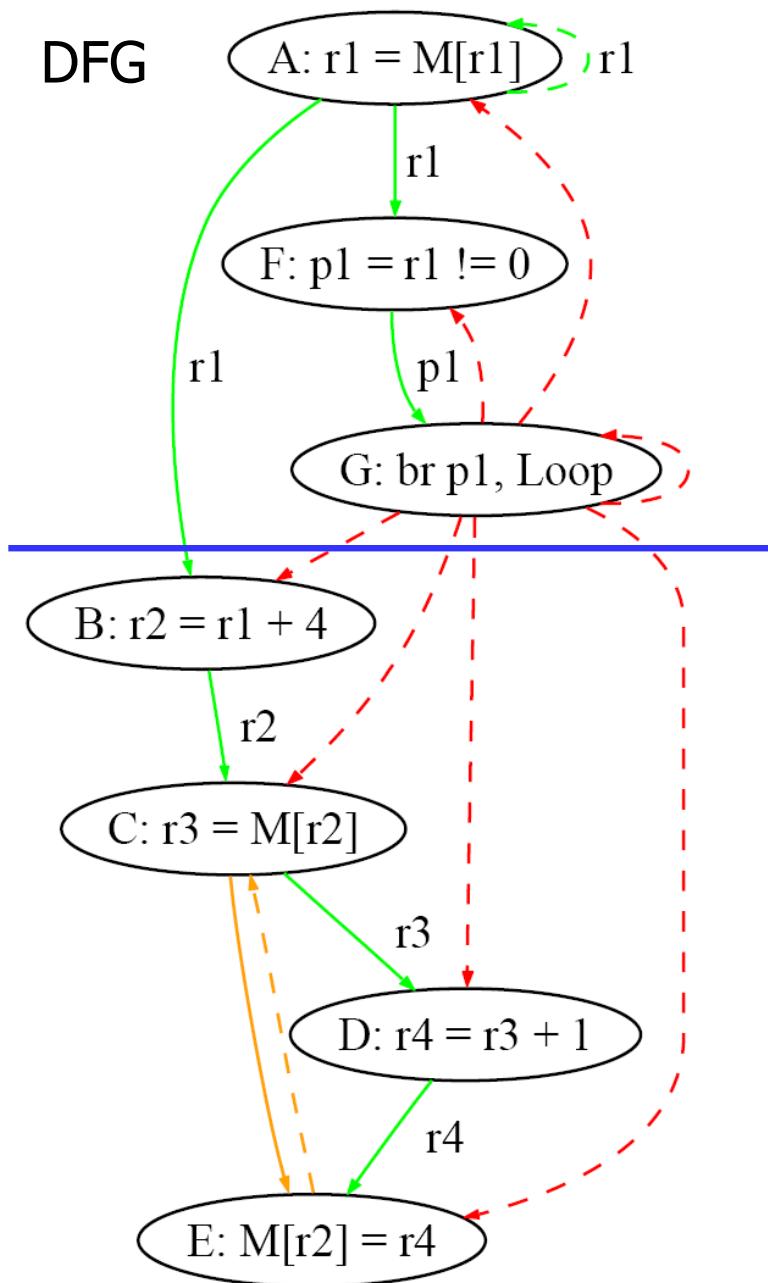
→ intra-iteration

→ loop-carried

█ communication queue

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

Implementing DSWP



L1:

SPAWN(Aux)
 A: $r1 = M[r1]$
 PRODUCE [1] = $r1$
 F: $p1 = r1 \neq 0$
 G: $\text{br } p1, L1$

Aux:

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

register

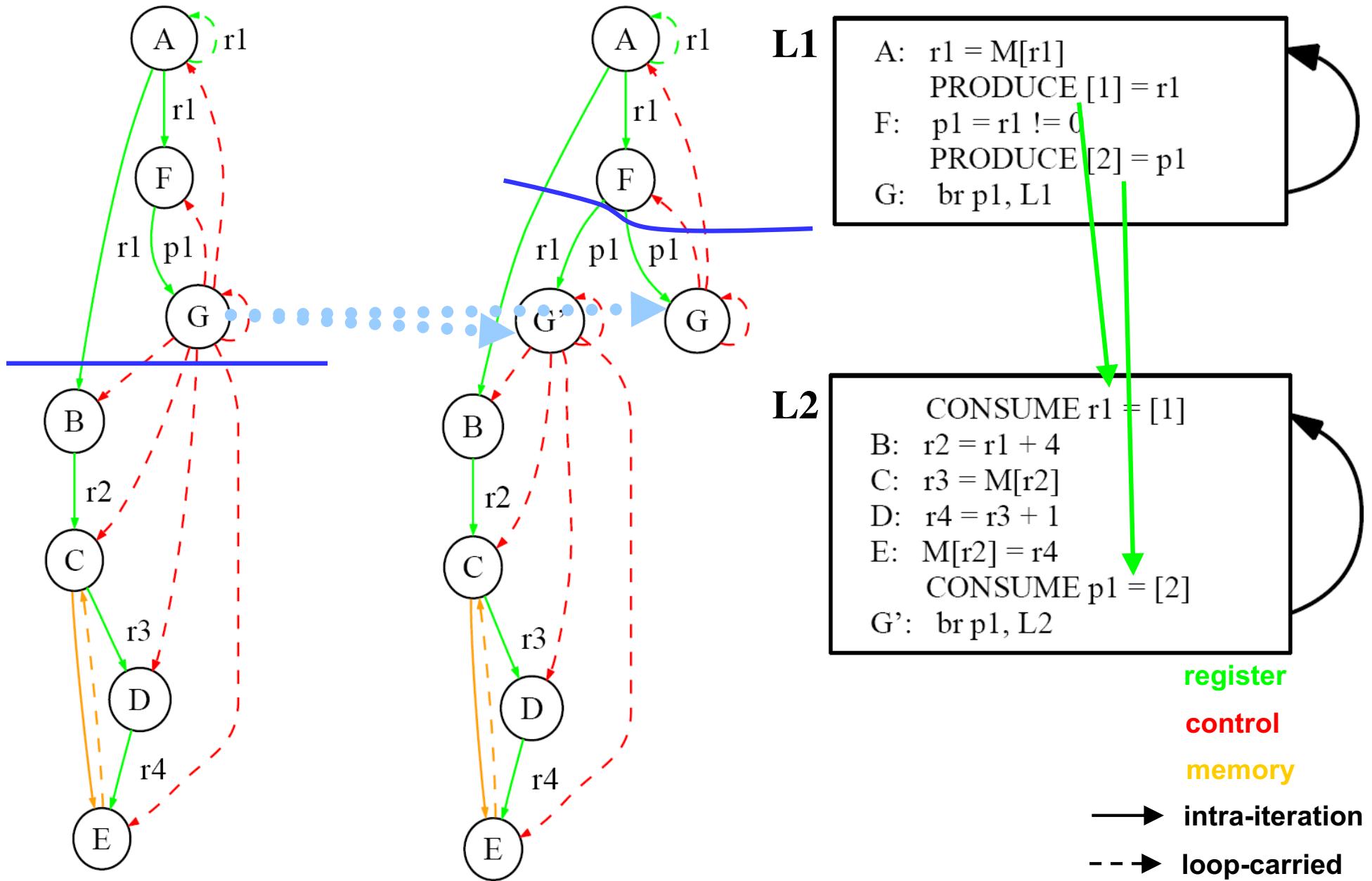
control

memory

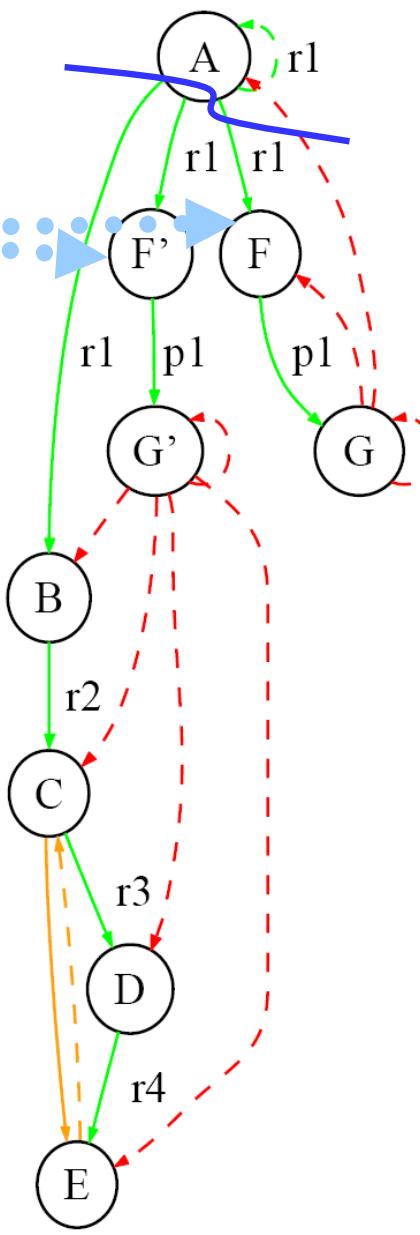
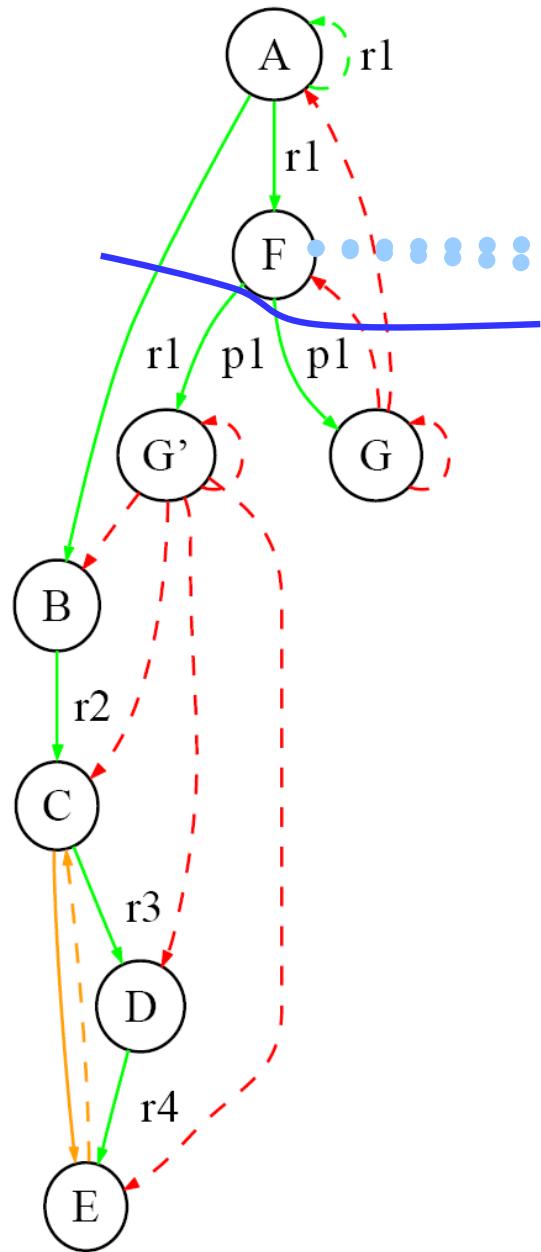
→ intra-iteration

→ loop-carried

Optimization: Node Splitting To Eliminate Cross Thread Control



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

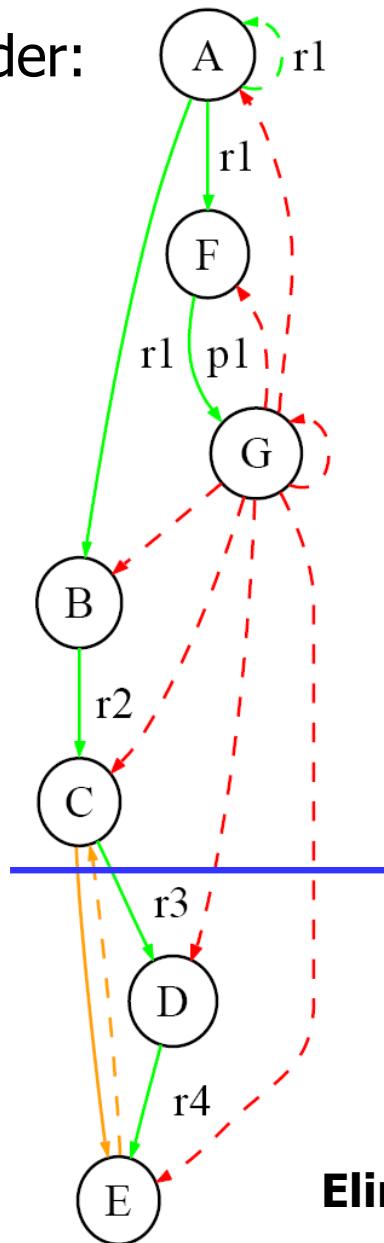
memory

→ intra-iteration

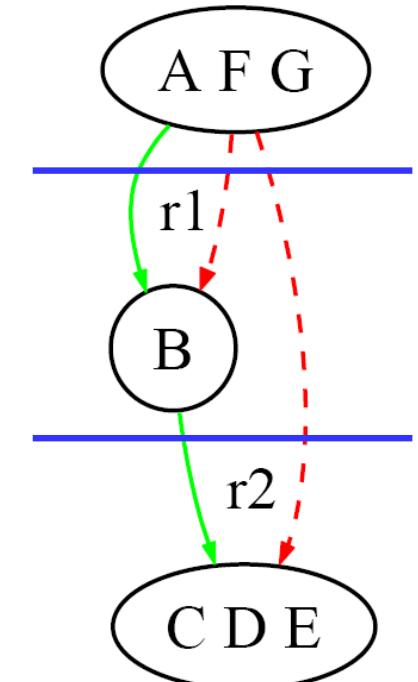
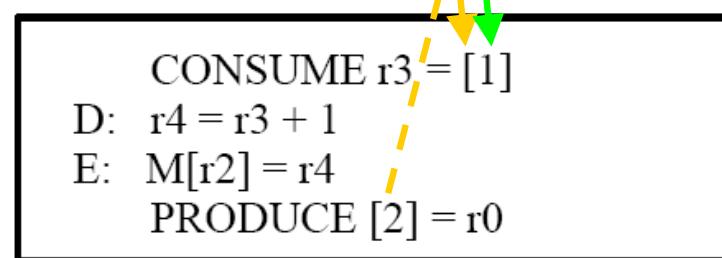
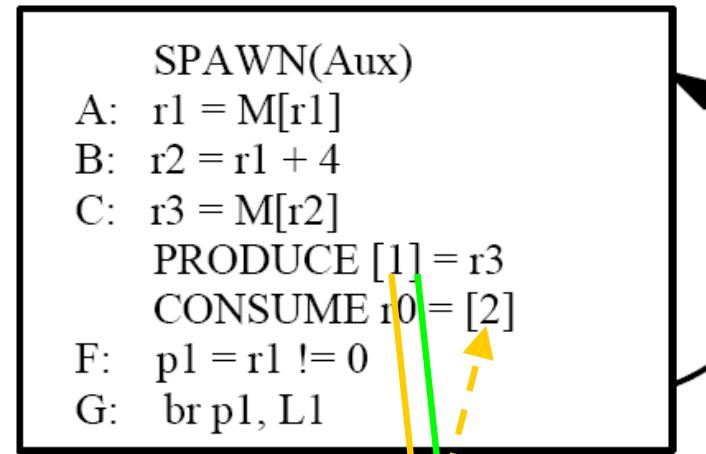
→ loop-carried

Constraint: Strongly Connected Components

Consider:



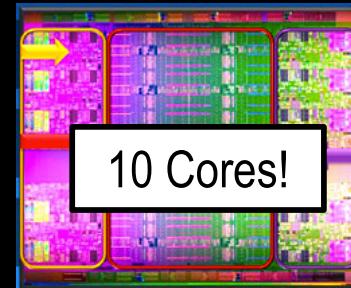
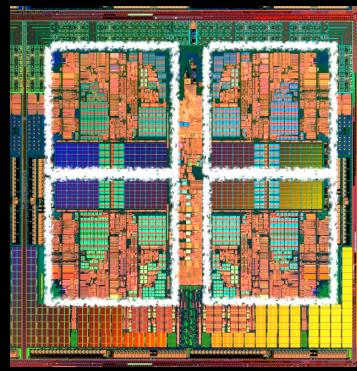
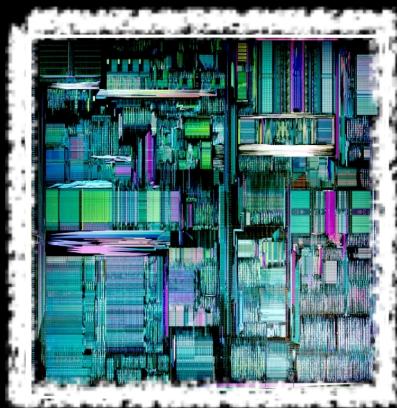
Solution: DAG_{SCC}



register
control
memory

→ intra-iteration
→ loop-carried

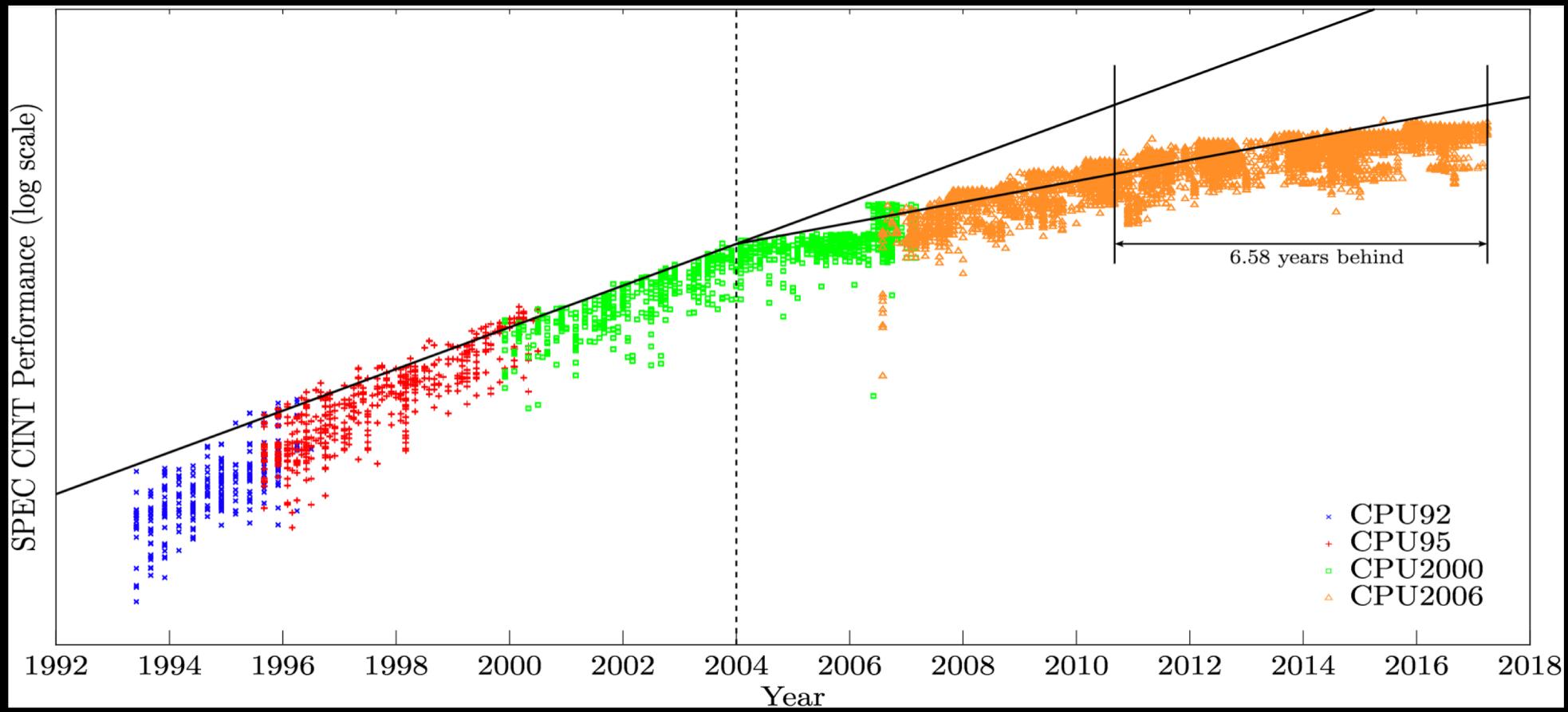
Eliminates pipelined/decoupled property



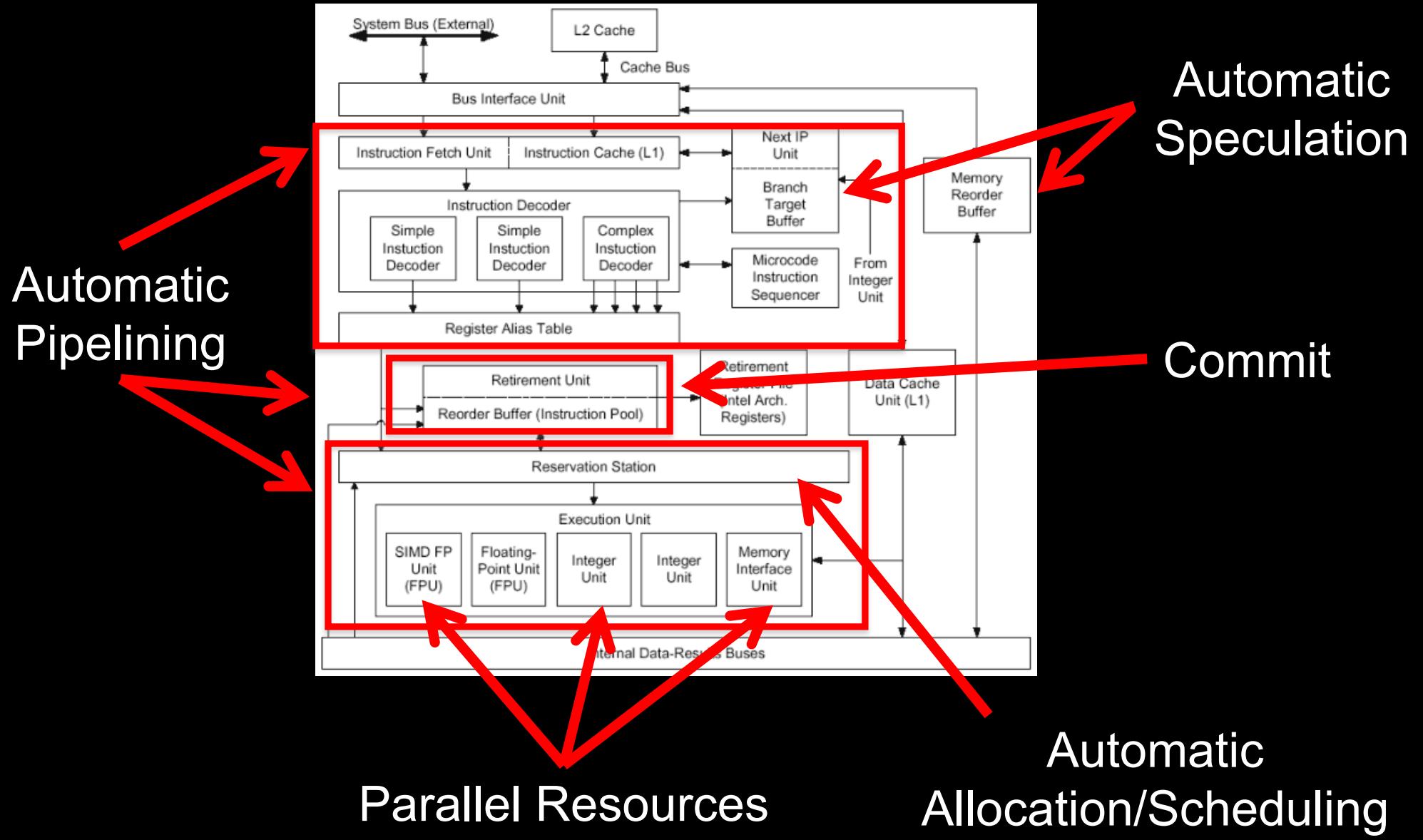
10-Core Intel Xeon
“Unparalleled Performance”

Era of DIY:

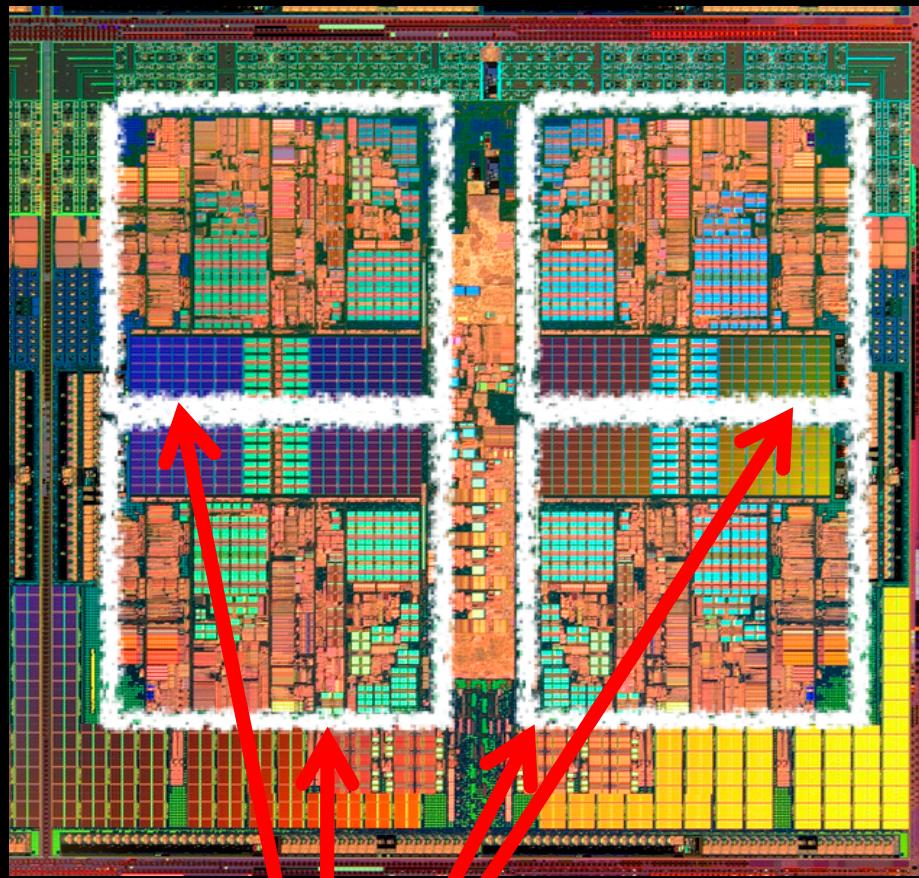
- Multicore
- Reconfigurable
- GPUs
- Clusters



P6 SUPERSCALAR ARCHITECTURE (CIRCA 1994)



MULTICORE ARCHITECTURE (CIRCA 2010)



Automatic
Pipelining

Parallel Resources

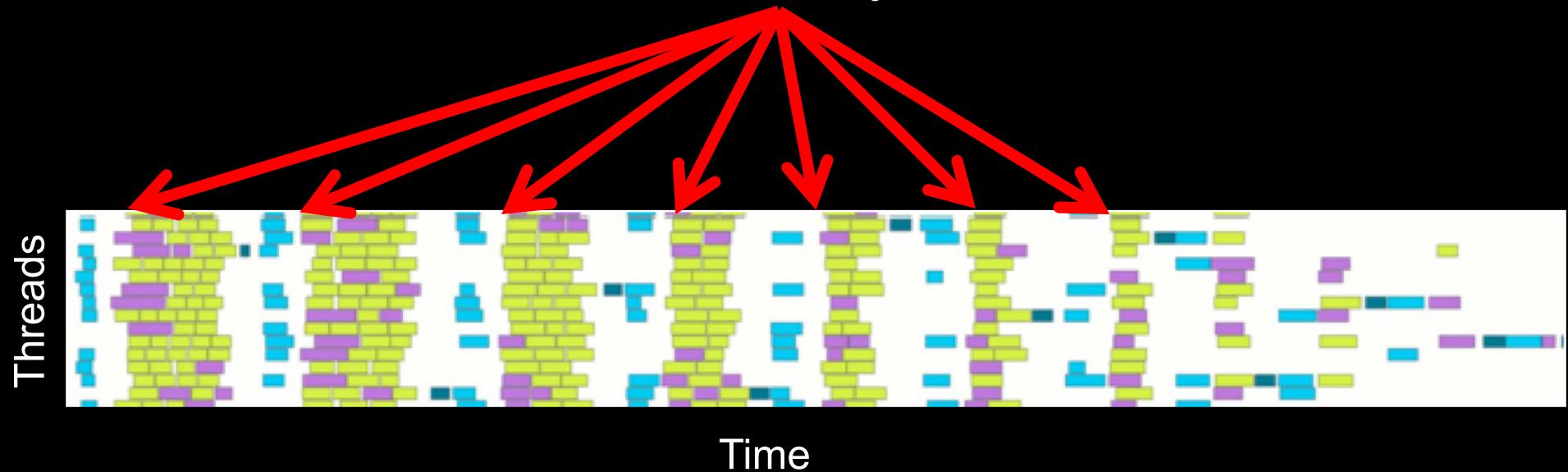
Automatic
Allocation/Scheduling

Automatic
Speculation

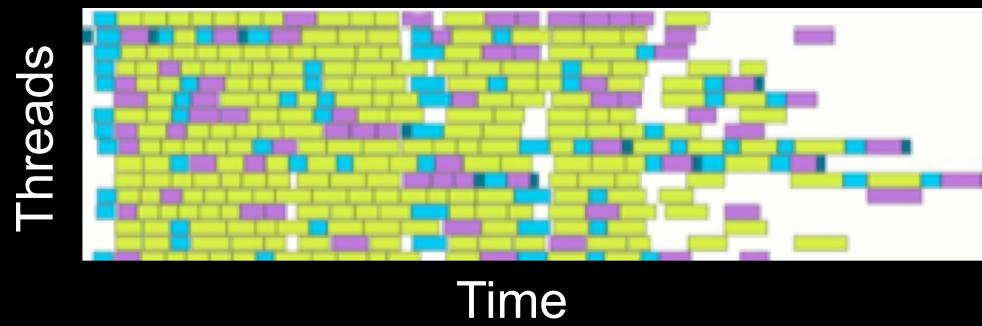
Commit

ABCPL	CORRELATE	GLU	Mentat	Paraphrase2	
ACE	CPS	GUARD	Legion	Paralation	pC++
ACT++	CRL	HaSL	Meta Chaos	Parallel-C++	SCHEDULE
Active messages	CSP	Haskell	Midway	Parallaxis	SciTL
Adl	Cthreads	HPC++	Millipede	ParC	POET
Adsmith	CUMULVS	JAVAR	CparPar	ParLib++	SDDA
ADDAP	DAGGER	HORUS	Mirage	ParLin	SHMEM
AFAPI	DAPPLE	HPC	MpC	Pannacs	SIMPLE
ALWAN	Data Parallel C	IMPACT	MOSIX	Parti	Sina
AM	DC++	ISIS	Modula-P	pC	SISAL
AMDC	DCE++	JAVAR	Modula-2*	pC++	distributed smalltalk
AppLeS	DDD	JADE	Multipol	PCN	SMI
Amoeba	DICE	Java RMI	MPI	PCP:	SONiC
ARTS	DIPC	javaPG	MPC++	PH	Split-C
Athapaskan-0b	DOLIB	JavaSpace	Munin	PEACE	SR
Aurora	DOME	JIDL	Nano-Threads	PCU	SThreads
Automap	DOSMOS	Joyce	NESL	PET	Strand
bb_threads	DRL	Khoros	NetClasses++	PETSc	SUIF
Blaze	DSM-Threads	Karma	Nexus	PENNY	Synergy
BSP	Ease	KOAN/Fortran-S	Nimrod	Phosphorus	Telegrphos
BlockComm	ECO	LAM	NOW	POET	SuperPascal
C*	Eiffel	Lilac	Objective Linda	Polaris	TCGMSG
"C* in C	Eilean	Linda	Occam	POOMA	Threads.h++
C**	Emerald	JADA	Omega	POOL-T	TreadMarks
CarlOS	EPL	WWWinda	OpenMP	PRESTO	TRAPPER
Cashmere	Excalibur	ISETL-Linda	Orca	P-RIO	uC++
C4	Express	ParLin	OOF90	Prospero	UNITY
CC++	Falcon	Eilean	P++	Proteus	UC
Chu	Filaments	P4-Linda	P3L	QPC++	V
Charlotte	FM	Glenda	p4-Linda	PVM	ViC*
Charm	FLASH	POSYBL	Pablo	PSI	Visifold V-NUS
Charm++	The FORCE	Objective-Linda	PADE	PSDM	VPE
Cid	Fork	LiPS	PADRE	Quake	Win32 threads
Cilk	Fortran-M	Locust	Panda	Quark	WinPar
CM-Fortran	FX	Lparx	Papers	Quick Threads	WWWinda
Converse	GA	Lucid	AFAPI	Sage++	XENOOPS
Code	GAMMA	Maisie	Para++	SCANDAL	XPC
COOL	Glenda	Manifold	Paradigm	SAM	Zounds
					ZPL

Parallel Library Calls

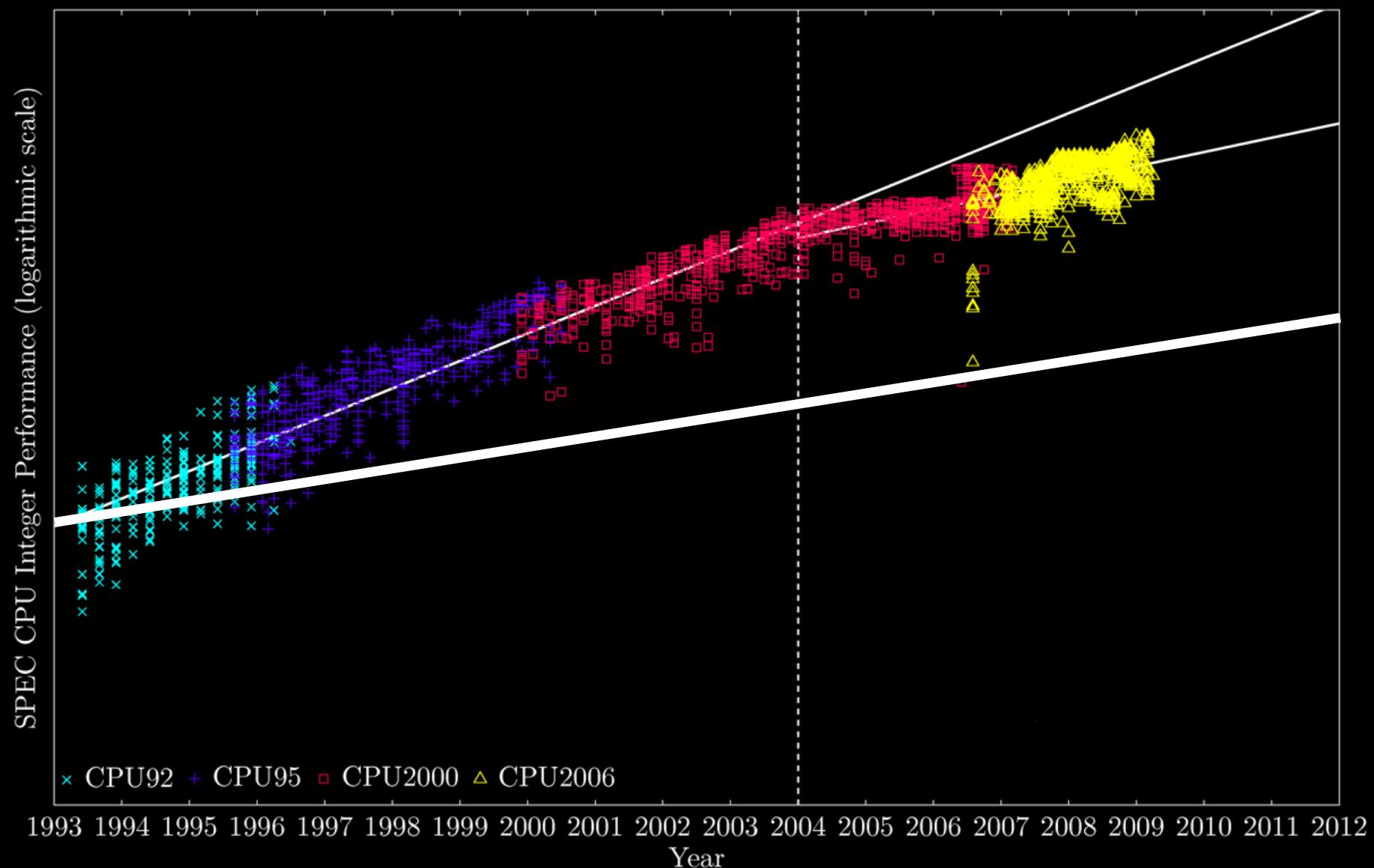


Realizable parallelism



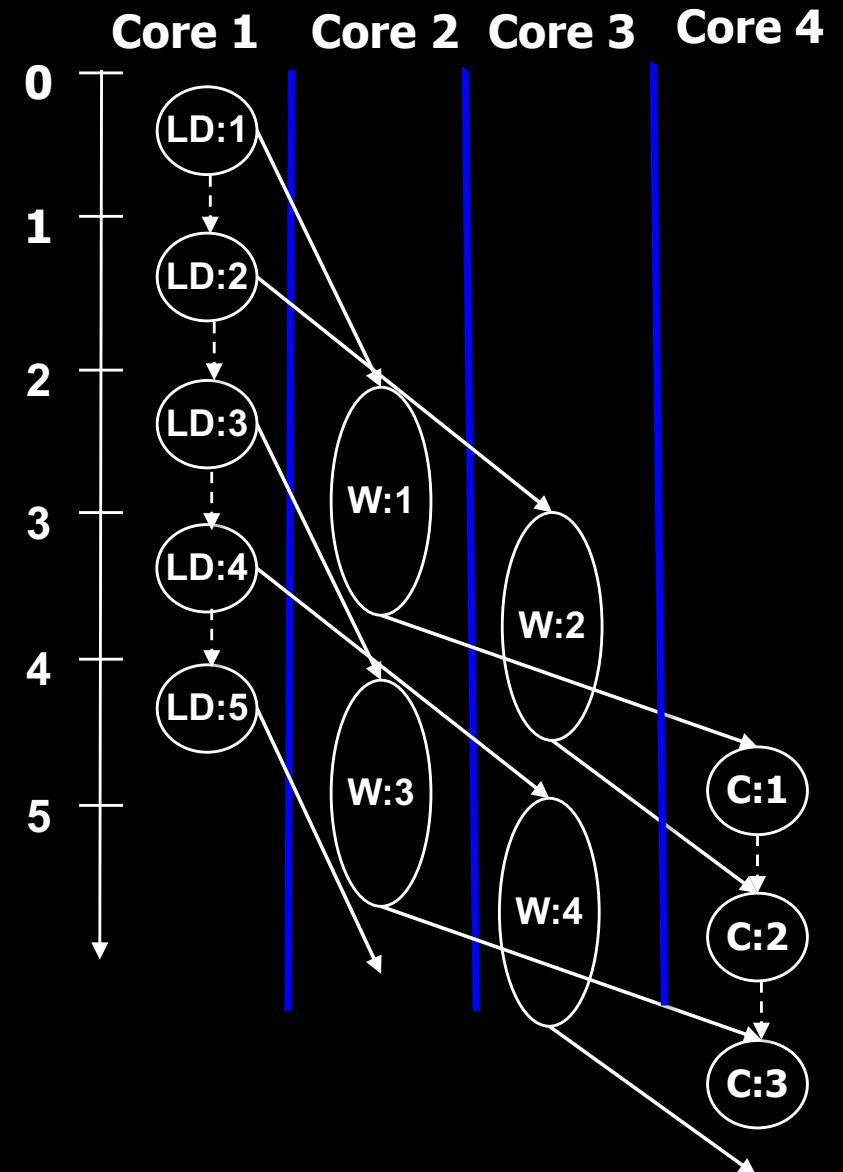
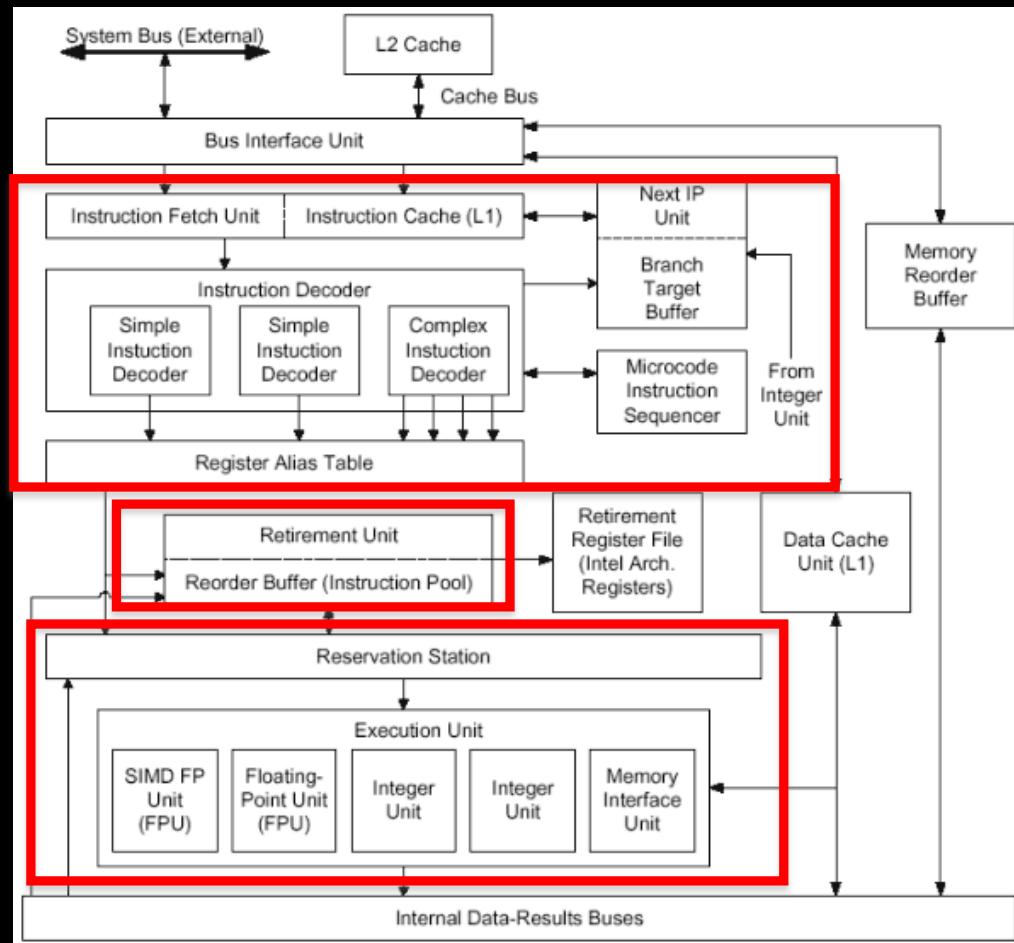
Credit: Jack Dongarra

“Compiler Advances Double Computing Power Every 18 Years!”
– Proebsting’s Law



Spec-PS-DSWP

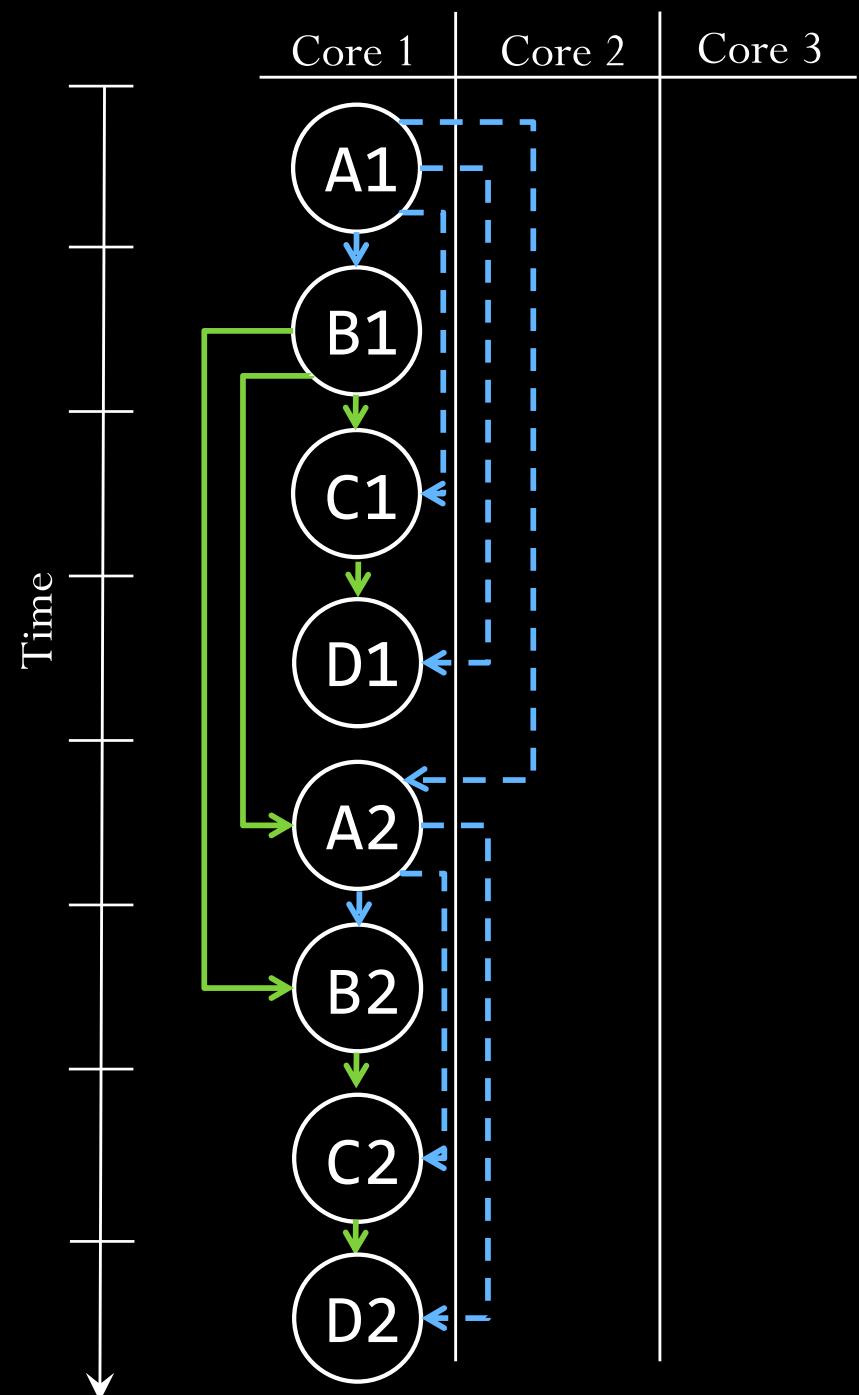
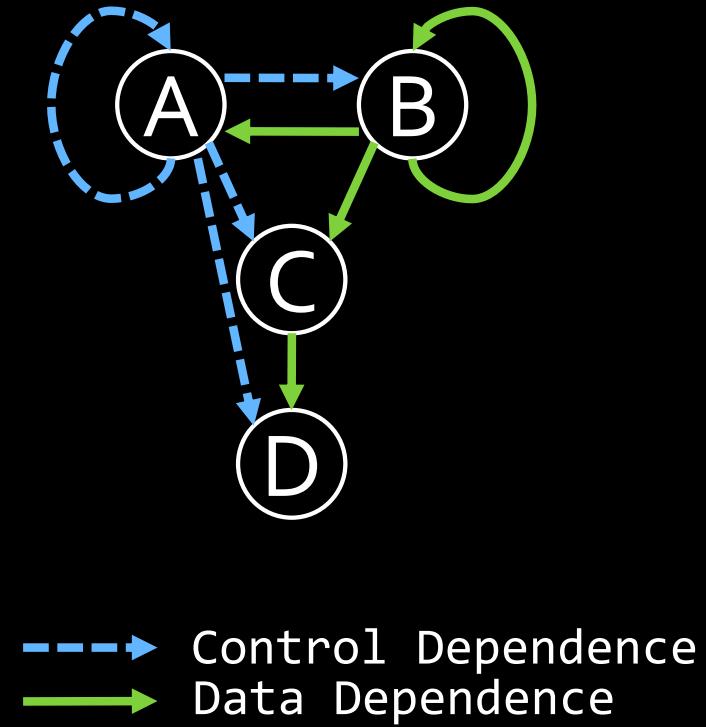
P6 SUPERSCALAR ARCHITECTURE



Example

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

Program Dependence Graph

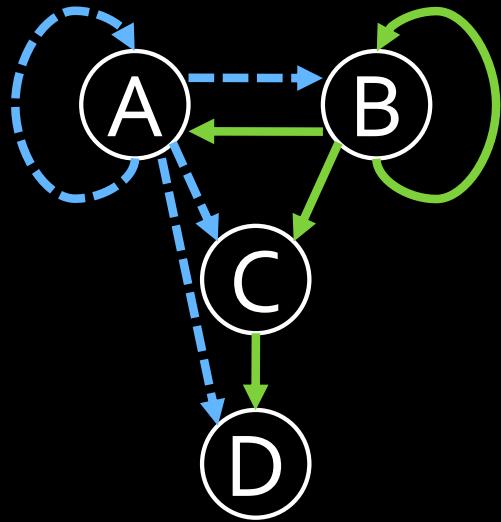


Spec-DOALL

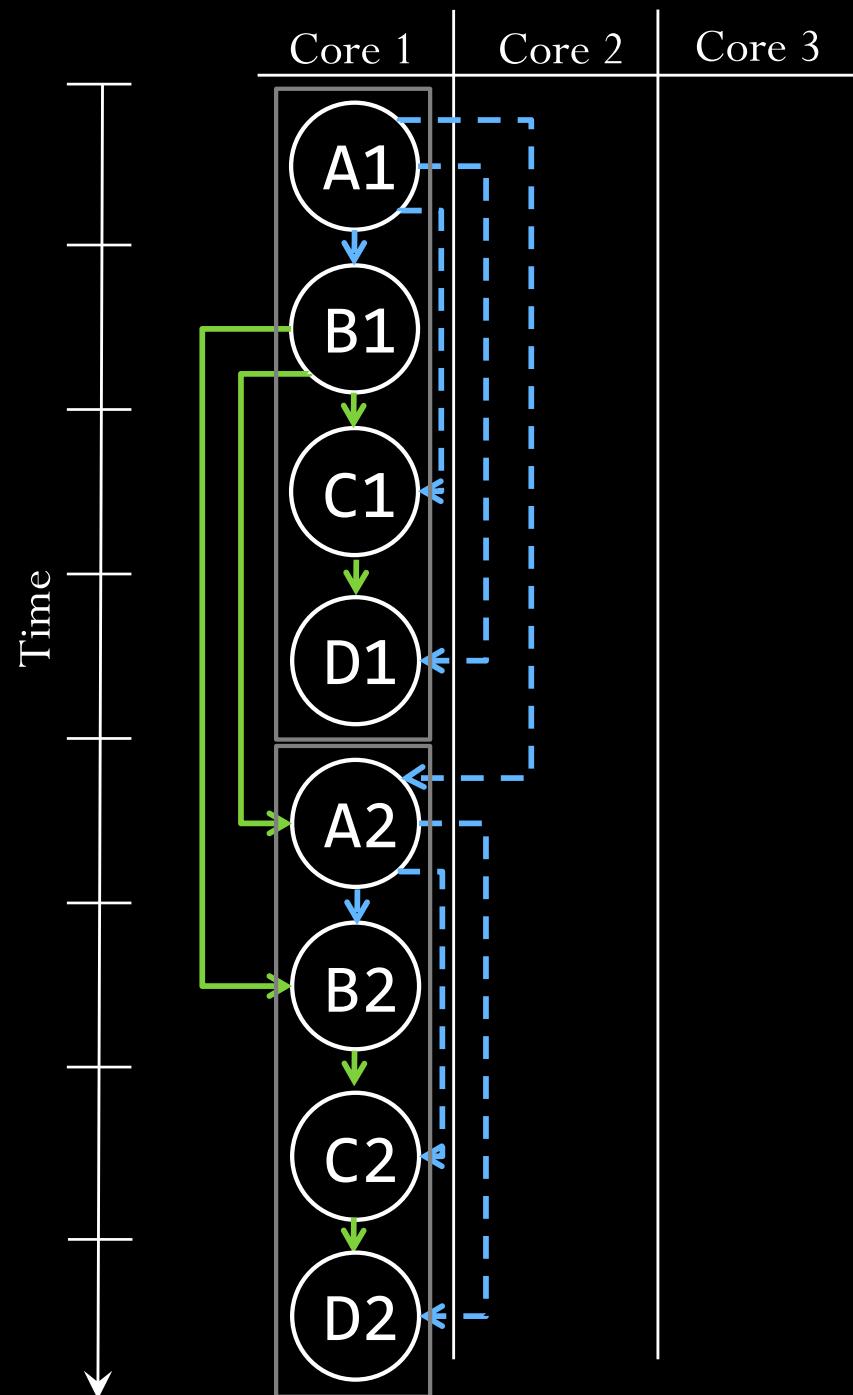
Example

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

Program Dependence Graph



—→ Control Dependence
—→ Data Dependence

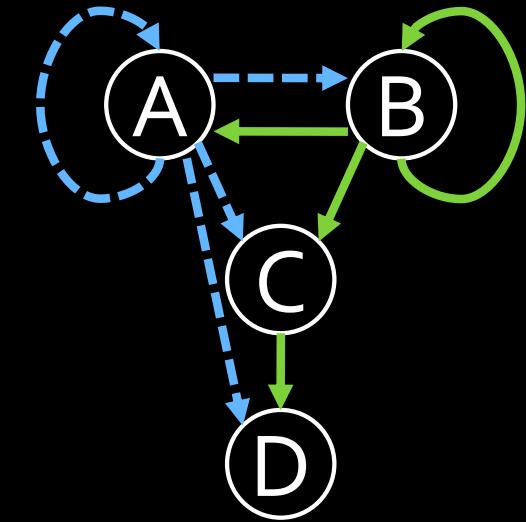


Spec-DOALL

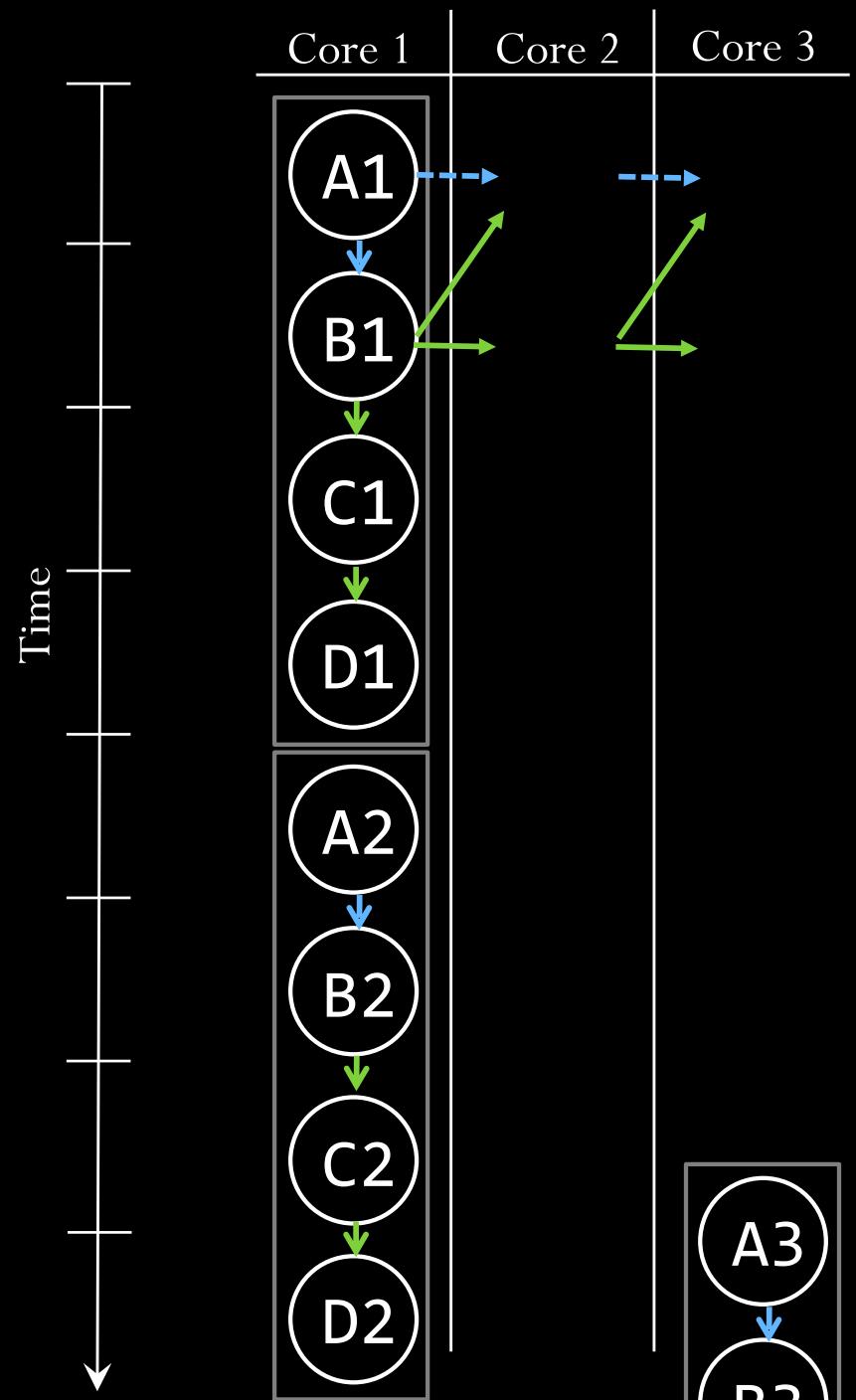
Example

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

Program Dependence Graph



—→ Control Dependence
—→ Data Dependence



Spec-DOALL

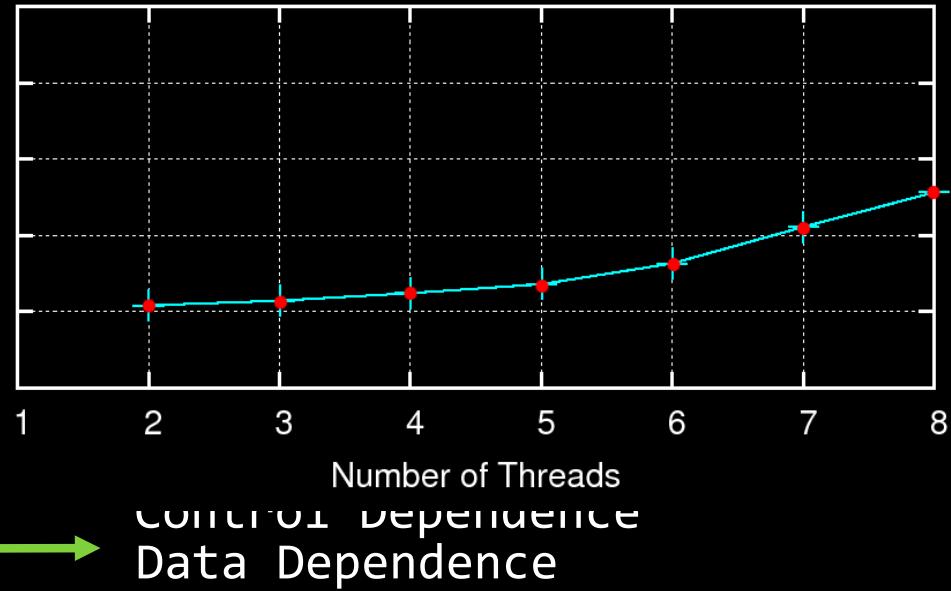
Example

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

Program Dependence Graph

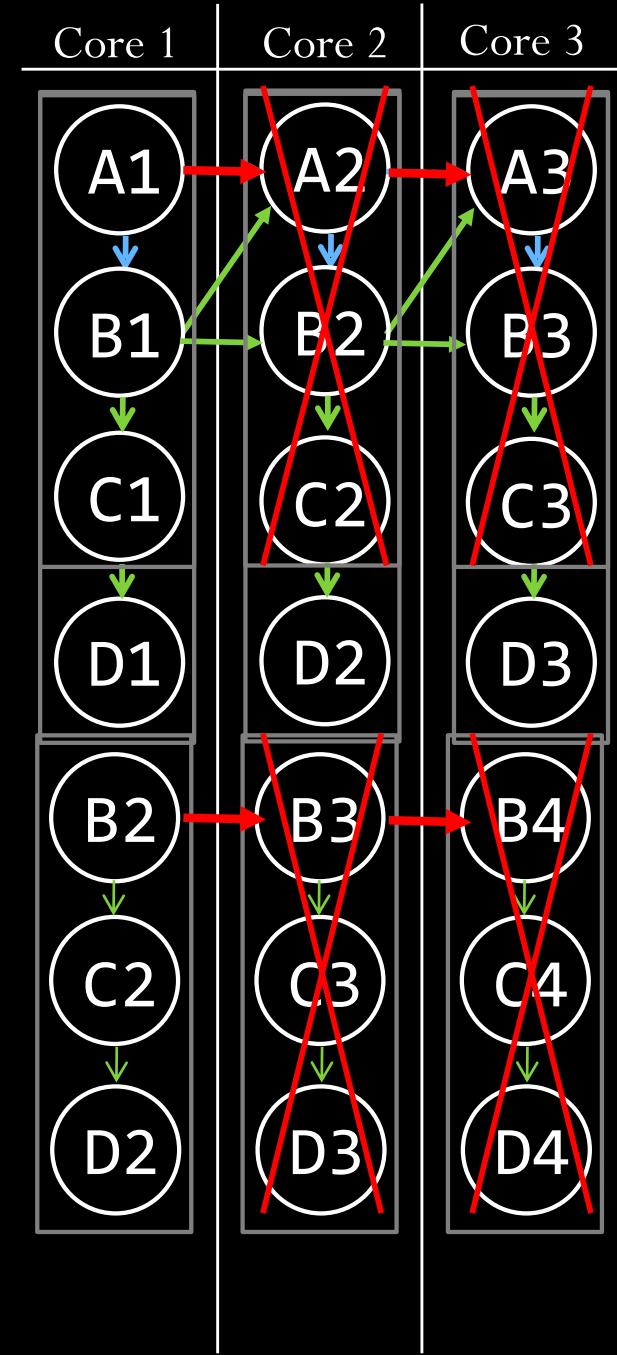


Slowdown



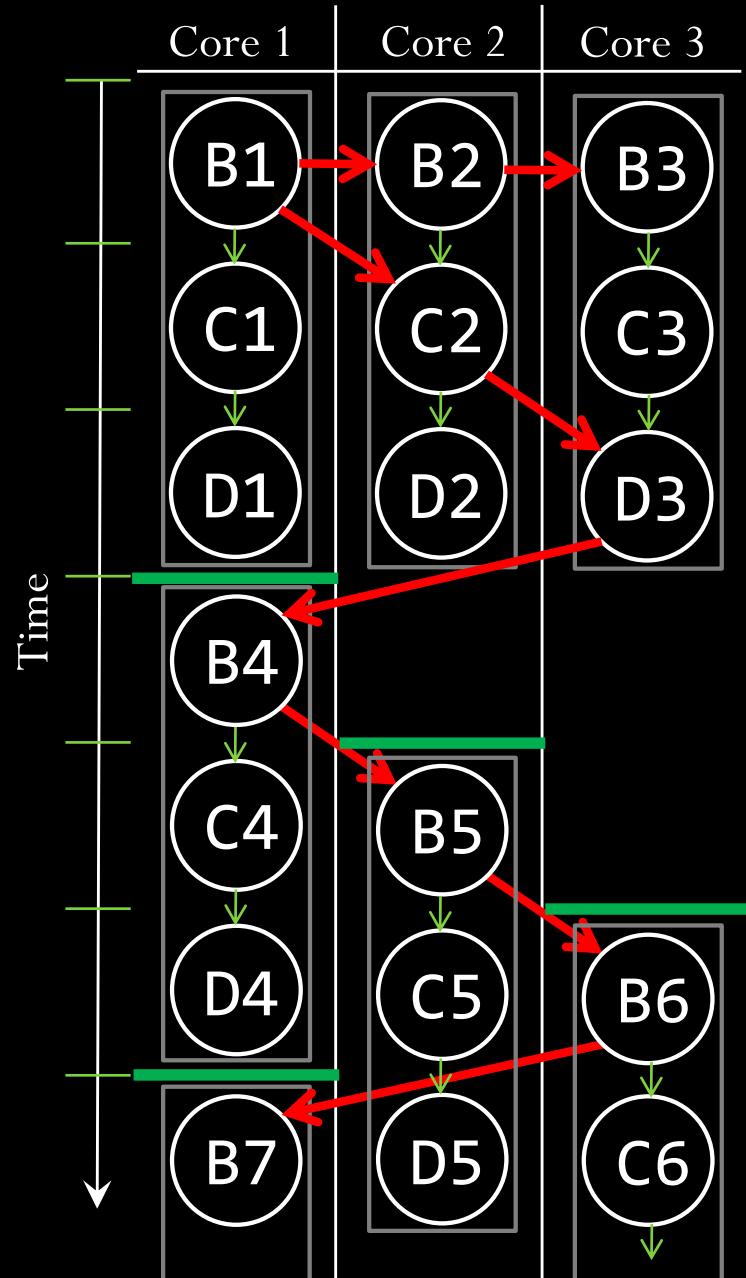
→ CONTROL DEPENDENCE
Data Dependence

Time ↓



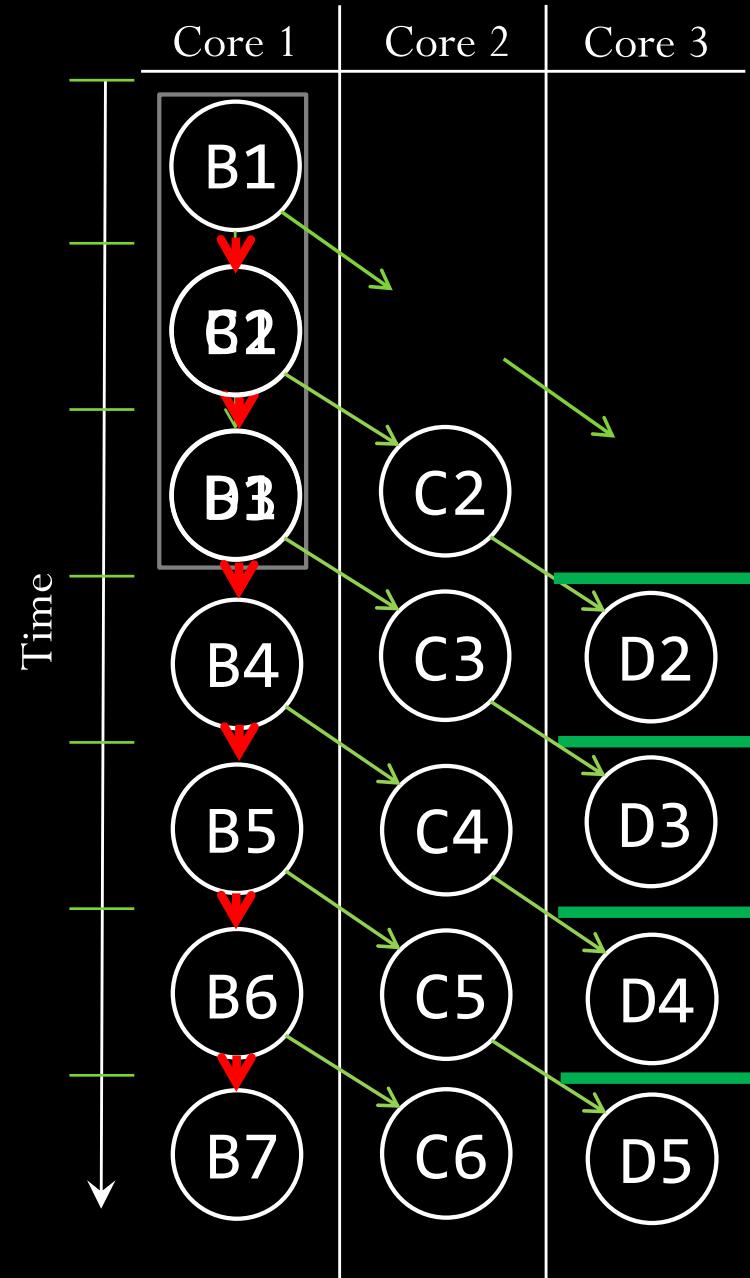
Spec-DOACROSS

Throughput: 1 iter/cycle



Spec-DSWP

Throughput: 1 iter/cycle



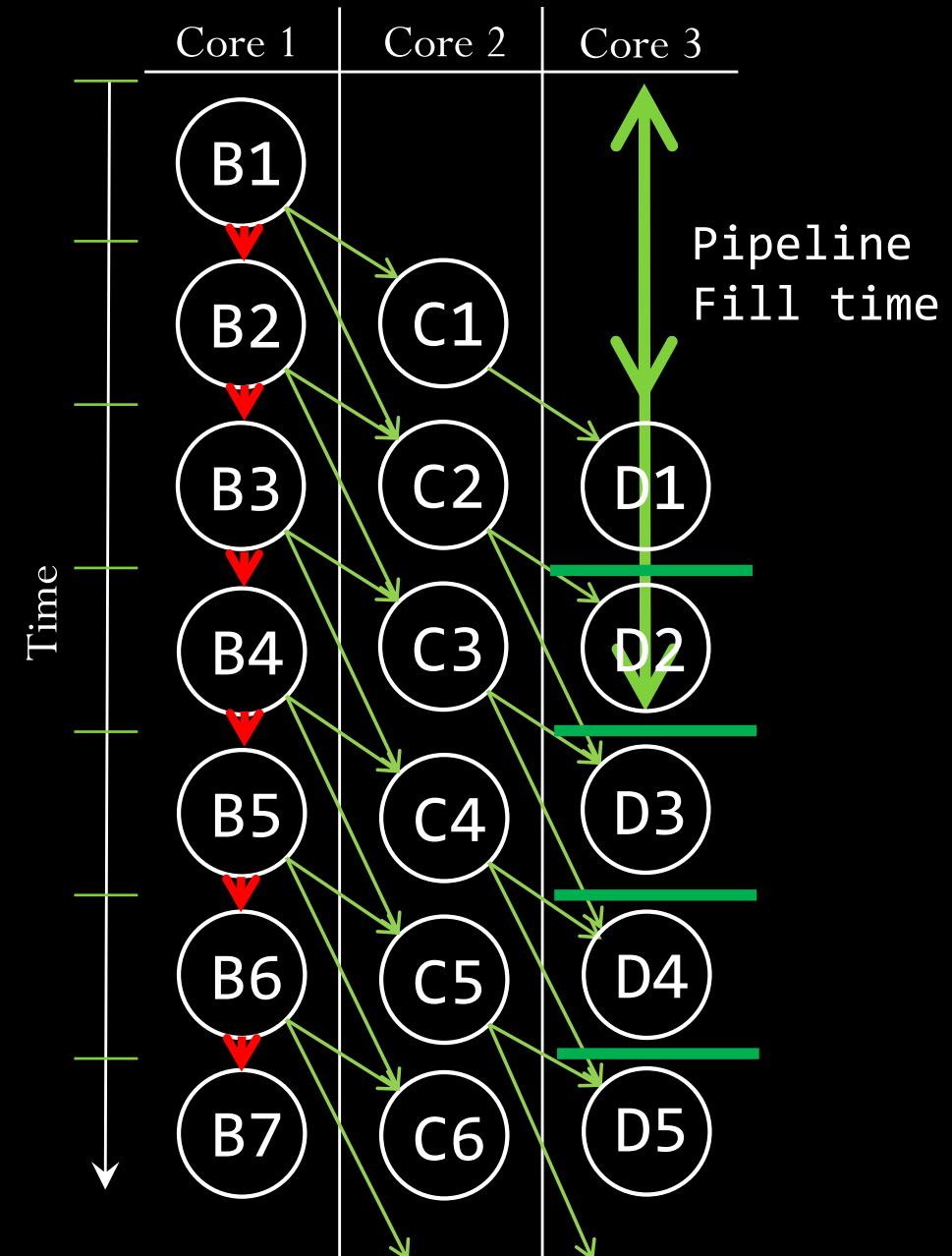
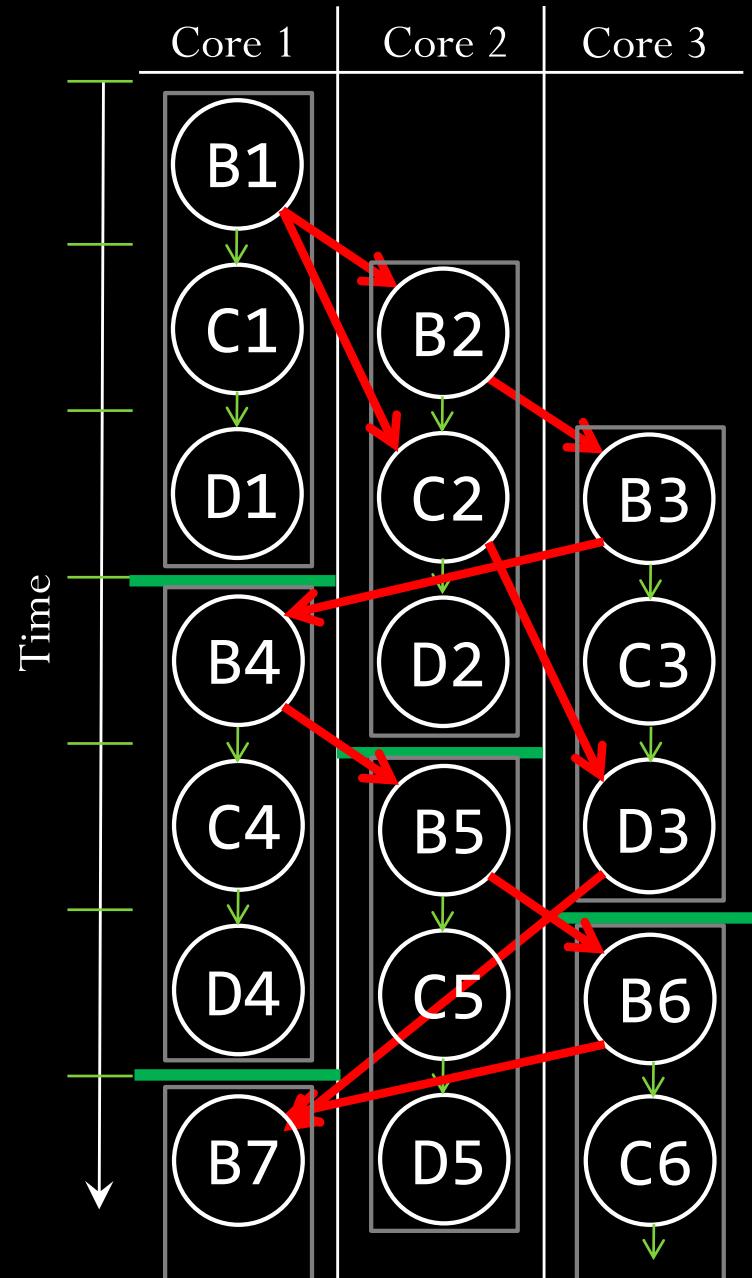
Comparison: Spec-DOACROSS and Spec-DSWP

Comm. Latency = 1: 1 iter/cycle

Comm. Latency = 2: **0.5 iter/cycle**

Comm. Latency = 1: 1 iter/cycle

Comm. Latency = 2: **1 iter/cycle**



Spec-DOACROSS vs. Spec-DSWP

[MICRO 2010]

Geomean of 11 benchmarks on the same cluster

