

# Lecture 1

## Introduction / Overview

EEC 171 Parallel Architectures

John Owens

UC Davis

# Credits

- © John Owens / UC Davis 2007–9.
- Thanks to many sources for slide material: Computer Organization and Design (Patterson & Hennessy) © 2005, Computer Architecture (Hennessy & Patterson) © 2007, Inside the Machine (Jon Stokes) © 2007, © Dan Connors / University of Colorado 2007, © Wen-Mei Hwu/David Kirk, University of Illinois 2007, © David Patterson / UCB 2003–6, © Mary Jane Irwin / Penn State 2005, © John Kubiawicz / UCB 2002, © Krste Asinovic/Arvind / MIT 2002, © Morgan Kaufmann Publishers 1998.

# Administrivia

# Teaching Staff

- Instructor

- John Owens | [jowens@ece](mailto:jowens@ece) | [www.ece/~jowens/](http://www.ece/~jowens/)

- Office hour: Kemper 3175, M 1–2

- TA

- Tracy Liu, [yliu@ucdavis](mailto:yliu@ucdavis)

- Office hour: Kemper 2243, T 2–3

# Electronic Resources

- [tinyurl.com/eec171-s09](http://tinyurl.com/eec171-s09) (points to SmartSite)
- Email class list (includes staff):
  - [eec171-s09@smartsite](mailto:eec171-s09@smartsite)
- Email teaching staff:
  - [eec171-s09-staff@listproc](mailto:eec171-s09-staff@listproc)
- Do not send mail to my personal account if you want a timely reply.

# Classroom Time

- Class is MW 2–4
- 2–3:30: “Lecture”
- 3:30–4: “Discussion”
- In reality: merged together
  - Discussion mostly will be used for lecture
  - Also for problem solving (TAs), quizzes, etc.
- Small class—let’s make it interactive!
- Administrative announcements: In middle of class

# Lecture Style

- Students prefer blackboard
- It's hard to show complex diagrams on a blackboard though
  - Plus I like slides better
  - I'll give you my notes—spend your time thinking not writing
- Might be using some Guided Notes
- Also will throw in some discussion questions
- Will use board when appropriate

# Textbook



- John L. Hennessy and David Patterson, “Computer Architecture: A Quantitative Approach”, 4th edition (Morgan Kaufmann), 2007
  - Don’t get the third edition



# Grading

- 3 segments to class:
  - Instruction level parallelism
  - Thread level parallelism
  - Data level parallelism
- Homework: 10% (3)
- Projects: 30% (3)
- Midterms: 15% each
- Final: 30% (cumulative)
- Goals:
  - Reduce homework dependence
  - Projects are important!

# Course Philosophy

- Third time I'm teaching this class
- Here's what I hope you'll get for any given technique in this class:
  - Understand *what* that technique is
  - Understand *why* that technique is important
  - Not understand (necessarily) *how* that technique is implemented

# Important Dates

- Midterms
  - M 27 April (concentrates on instruction-level parallelism)
  - W 27 May (concentrates on thread-level parallelism)
  - TA will administer exams
- Final
  - W 10 June (6–8p)
  - Cumulative
- Exams are open-book, open-note
- Makeups on midterms or final are oral only

# Homework Turn-In

- Homework goes in 2131 Kemper
- Homework is due at noon
- Written homework must be done individually
- Homework and exam solutions will be handed out in class
- We'll try to make solutions ASAP after due date
- Please do not pass these solutions on to other students
- Use of online solutions to homework/projects is cheating

# Project Turnin

- Projects will be turned in electronically (SmartSite)
- Project deliverable will be a *writeup*
  - Ability to communicate is important!
  - Writeups will be short (1 page)
  - PDF
- Projects are individual

# Homework 0

- Due 5 pm Tuesday
- Please link a photo to your SmartSite profile
- Give yourself a big head

# What is cheating?

- Cheating is claiming credit for work that is not your own.
- Cheating is disobeying or subverting instructions of the instructional staff.
  - Homework deadlines, online solutions, etc.
- It is OK to work in (small) groups on homework.
  - All work you turn in must be 100% yours, and you must be able to explain all of it.
  - Give proper credit if credit is due.

# Things You Should Do

- *Ask questions!*
  - Especially when things aren't clear
- Give feedback!
  - Email or face-to-face
  - Tell me what I'm doing poorly
  - Tell me what I'm doing well
  - Tell the TA too
- Start projects early



# Things You Shouldn't Do

- Cheat
- Skip class (I'll know when you're not there!)
- Be late for class
- Read the paper in class
- Allow your cell phone to ring in class
- Ask OH questions without preparing
  - Make sure you do the reading!
  - Identifying what you have trouble with helps me

# Getting the Grade You Want

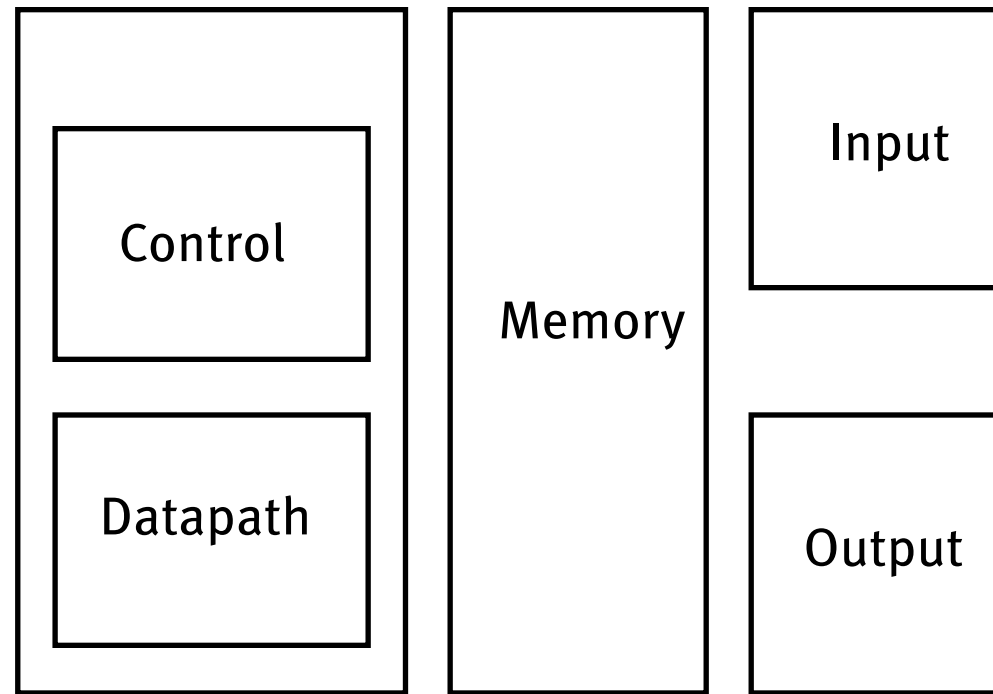
- Come to class!
- Ask questions in class when you don't understand
- Come to office hours (mine and Tracy's)
- Start the hw and projects early
- Use the projects as a vehicle for learning
- Understand the course material

# My Expectations

- This class was hard
- I learned a lot
- The work I did in this class was worthwhile
- The instructor was fair
- The instructor was effective
- The instructor cared about my learning

# Review

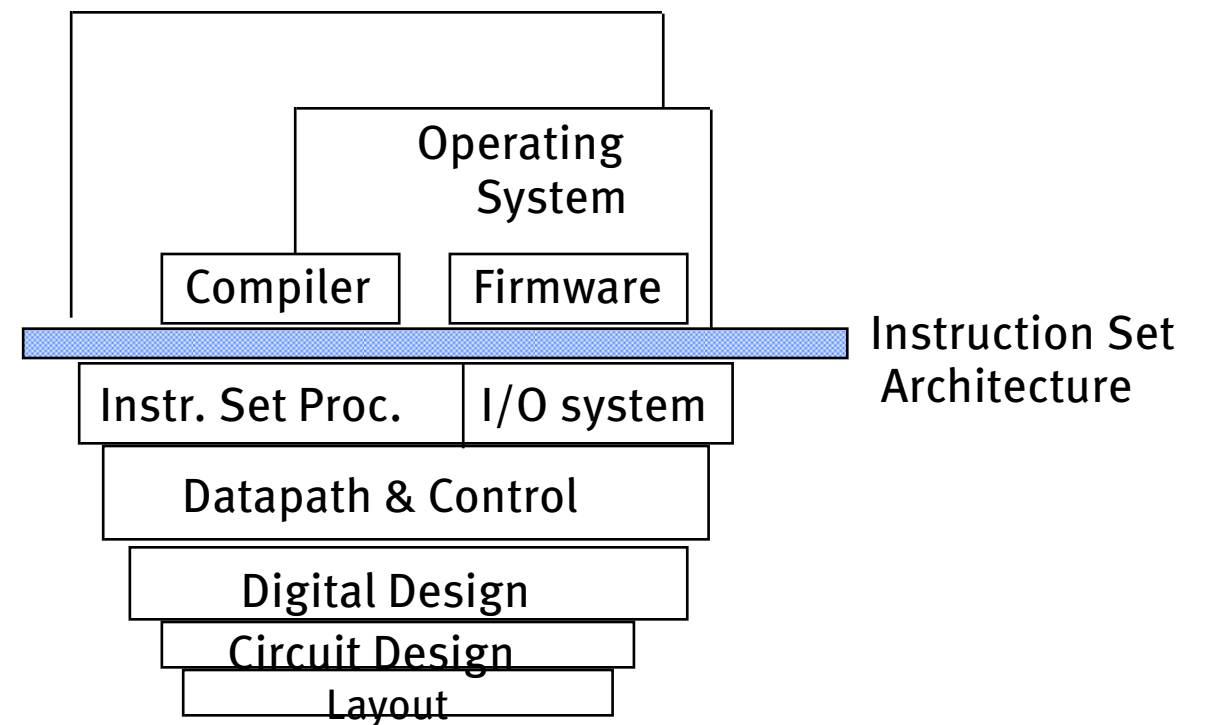
# The Big Picture



- Since 1946 all computers have had 5 components

# What is “Computer Architecture”?

- Coordination of many levels of abstraction
- Under a rapidly changing set of forces
- Design, Measurement, and Evaluation

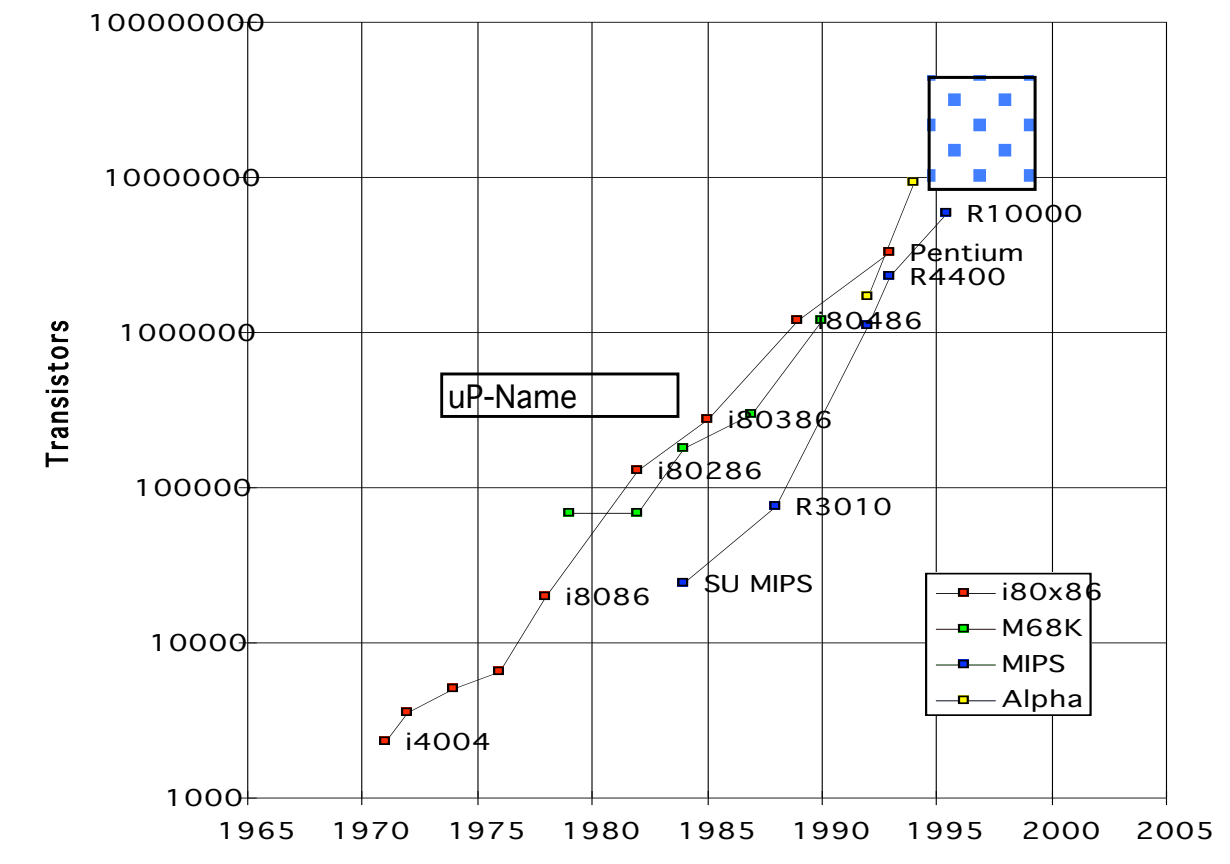


# Technology

## DRAM chip capacity

<u>Year</u>	<u>Size</u>
1980	64 Kb
1983	256 Kb
1986	1 Mb
1989	4 Mb
1992	16 Mb
1996	64 Mb
1999	256 Mb
2002	1 Gb
2005	4 Gb

## Microprocessor Logic Density



- In ~1985 the single-chip processor (32-bit) and the single-board computer emerged
- => workstations, personal computers, multiprocessors have been riding this wave since

# Technology rates of change

- Processor
  - logic capacity: about 30% per year
  - clock rate: about 20% per year
- Memory
  - DRAM capacity: about 60% per year (4x every 3 years)
  - Memory speed: about 10% per year
  - Cost per bit: improves about 25% per year
- Disk
  - capacity: about 60% per year
  - Total use of data: 100% per 9 months!
- Network Bandwidth increasing more than 100% per year!



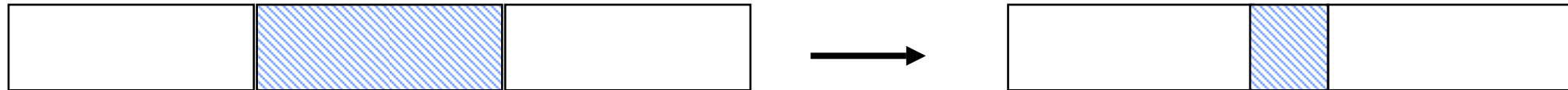
# The Performance Equation

- Time = Clock Speed \* CPI \* Instruction Count
  - = seconds/cycle \* cycles/instr \* instrs/program
  - => seconds/program
- “The only reliable measure of computer performance is time.”

# Amdahl's Law

- Speedup due to enhancement  $E$ :

$$\text{Speedup}(E) = \frac{\text{Execution Time without } E}{\text{Execution Time with } E} = \frac{\text{Performance with } E}{\text{Performance without } E}$$



- Suppose that enhancement  $E$  accelerates a fraction  $F$  of the task by a factor  $S$  and the remainder of the task is unaffected:

$$\text{Execution time (with } E) = ((1 - F) + F/S) \cdot \text{Execution time (without } E)$$

$$\text{Speedup (with } E) = \frac{1}{(1 - F) + F/S}$$

- Design Principle: Make the common case fast!

# Amdahl's Law example

- New CPU 10X faster
- I/O bound server, so 60% time waiting for I/O

$$\begin{aligned}\text{Speedup}_{\text{overall}} &= \frac{1}{(1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}}} \\ &= \frac{1}{(1 - 0.4) + \frac{0.4}{10}} = \frac{1}{0.64} = 1.56\end{aligned}$$

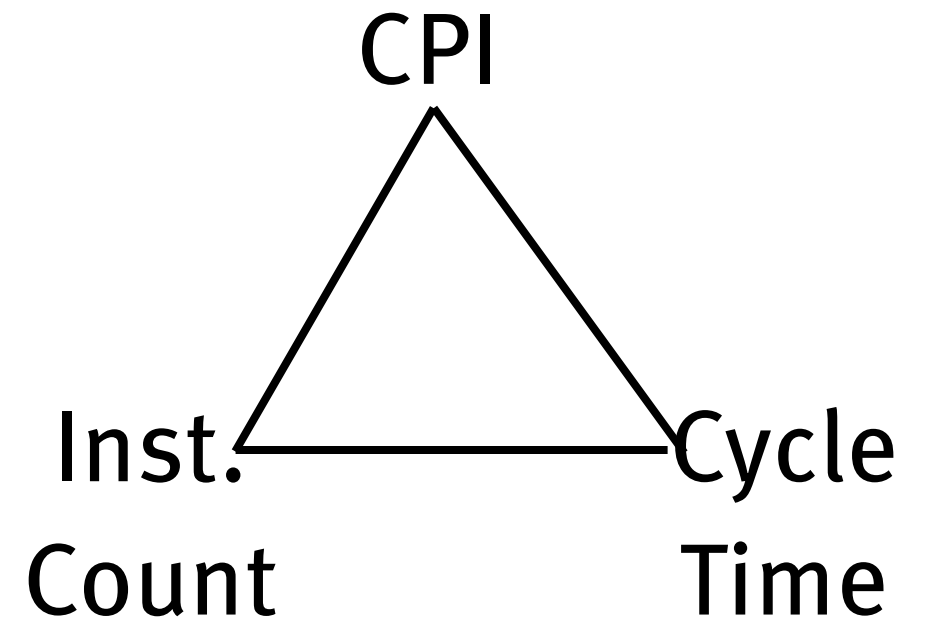
- Apparently, it's human nature to be attracted by 10X faster, vs. keeping in perspective it's just 1.6X faster

# Basis of Evaluation

Pros		Cons
<ul style="list-style-type: none"><li>• representative</li></ul>	Actual Target Workload	<ul style="list-style-type: none"><li>• very specific</li><li>• non-portable</li><li>• difficult to run, or measure</li><li>• hard to identify cause</li></ul>
<ul style="list-style-type: none"><li>• portable</li><li>• widely used</li><li>• improvements useful in reality</li></ul>	Full Application Benchmarks	<ul style="list-style-type: none"><li>• less representative</li></ul>
<ul style="list-style-type: none"><li>• easy to run, early in design cycle</li></ul>	Small “kernel” benchmarks	<ul style="list-style-type: none"><li>• easy to “fool”</li></ul>
<ul style="list-style-type: none"><li>• identify peak capability and potential bottlenecks</li></ul>	Microbenchmarks	<ul style="list-style-type: none"><li>• “peak” may be a long way from application performance</li></ul>

# Evaluating Instruction Sets

- Design-time Metrics:
  - Can it be implemented, in how long, at what cost?
  - Can it be programmed? Ease of compilation?
- Static Metrics:
  - How many bytes does the program occupy in memory?
- Dynamic Metrics:
  - How many instructions are executed?
  - How many bytes does the processor fetch to execute the program?
  - How many clocks are required per instruction?
  - How “lean” a clock is practical?
- Best Metric: Time to execute the program!



# MIPS Instruction Set

- 32-bit fixed format inst (3 formats)
- 32 32-bit GPR (R0 contains zero); 32 FP registers (and HI LO)
- partitioned by software convention
- 3-address, reg-reg arithmetic instr.
- Single address mode for load/store: base+displacement
  - no indirection, scaled
- 16-bit immediate plus LUI
- Simple branch conditions
  - compare against zero or two registers for =, ≠
  - no integer condition codes
- Delayed branch
  - execute instruction after a branch (or jump) even if the branch is taken
  - Compiler can fill branch delay slot ~50% of the time

# RISC Philosophy

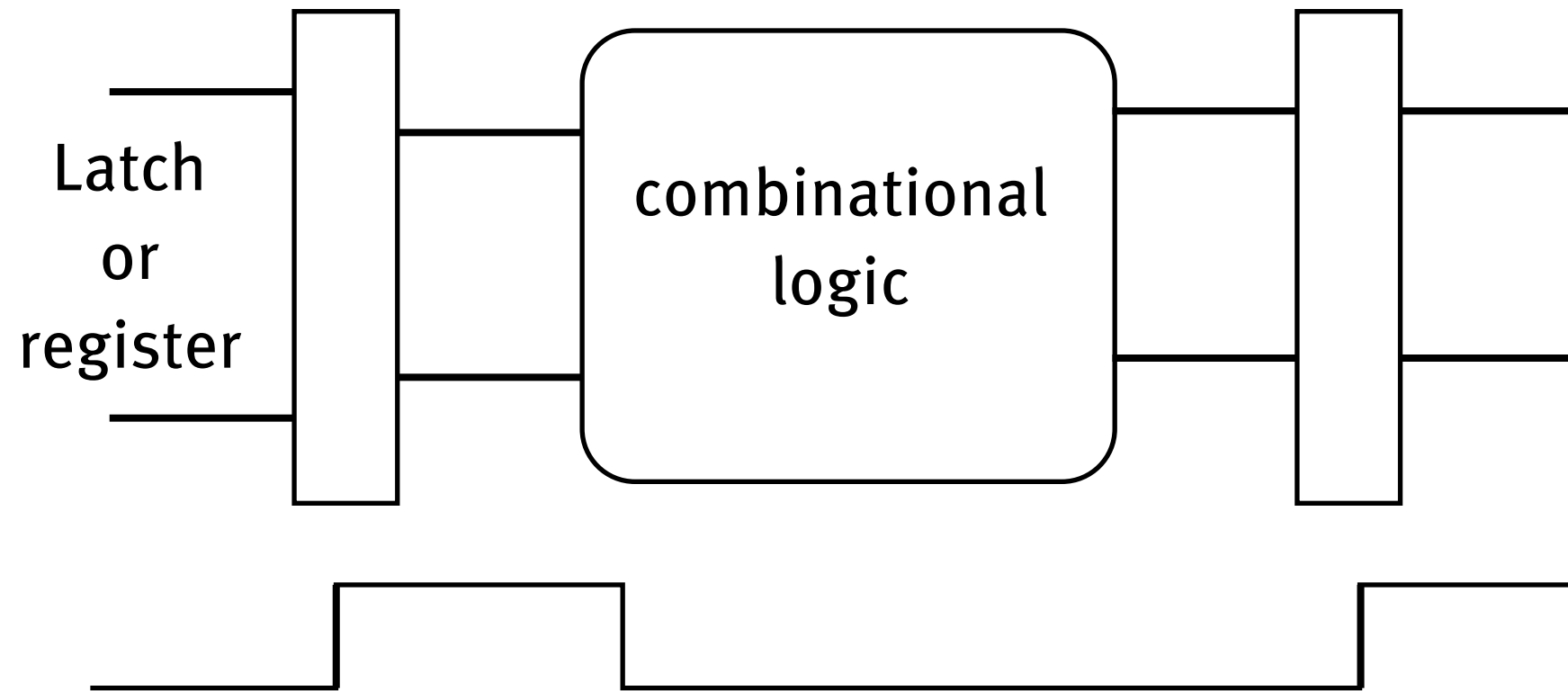
- Instructions all same size
- Small number of opcodes (small opcode space)
- Opcode in same place for every instruction
- Simple memory addressing
- Instructions that manipulate data don't manipulate memory, and vice versa
- Minimize memory references by providing ample registers

# Computer Arithmetic

- Bits have no inherent meaning: operations determine whether really ASCII characters, integers, floating point numbers
- 2's complement
- Hardware algorithms for arithmetic:
  - Carry lookahead/carry save addition (parallelism!)
- Floating point

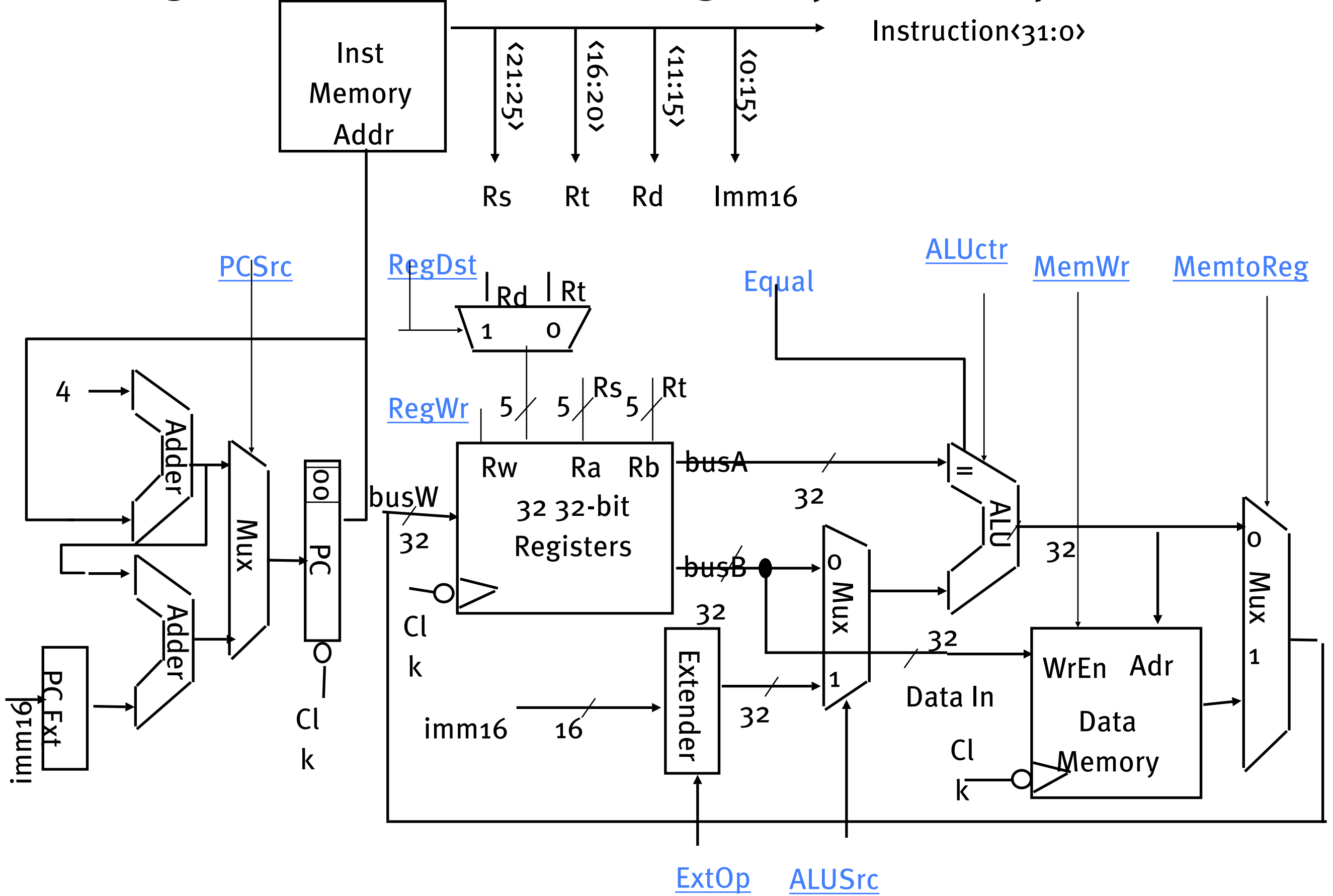


# What's a Clock Cycle?

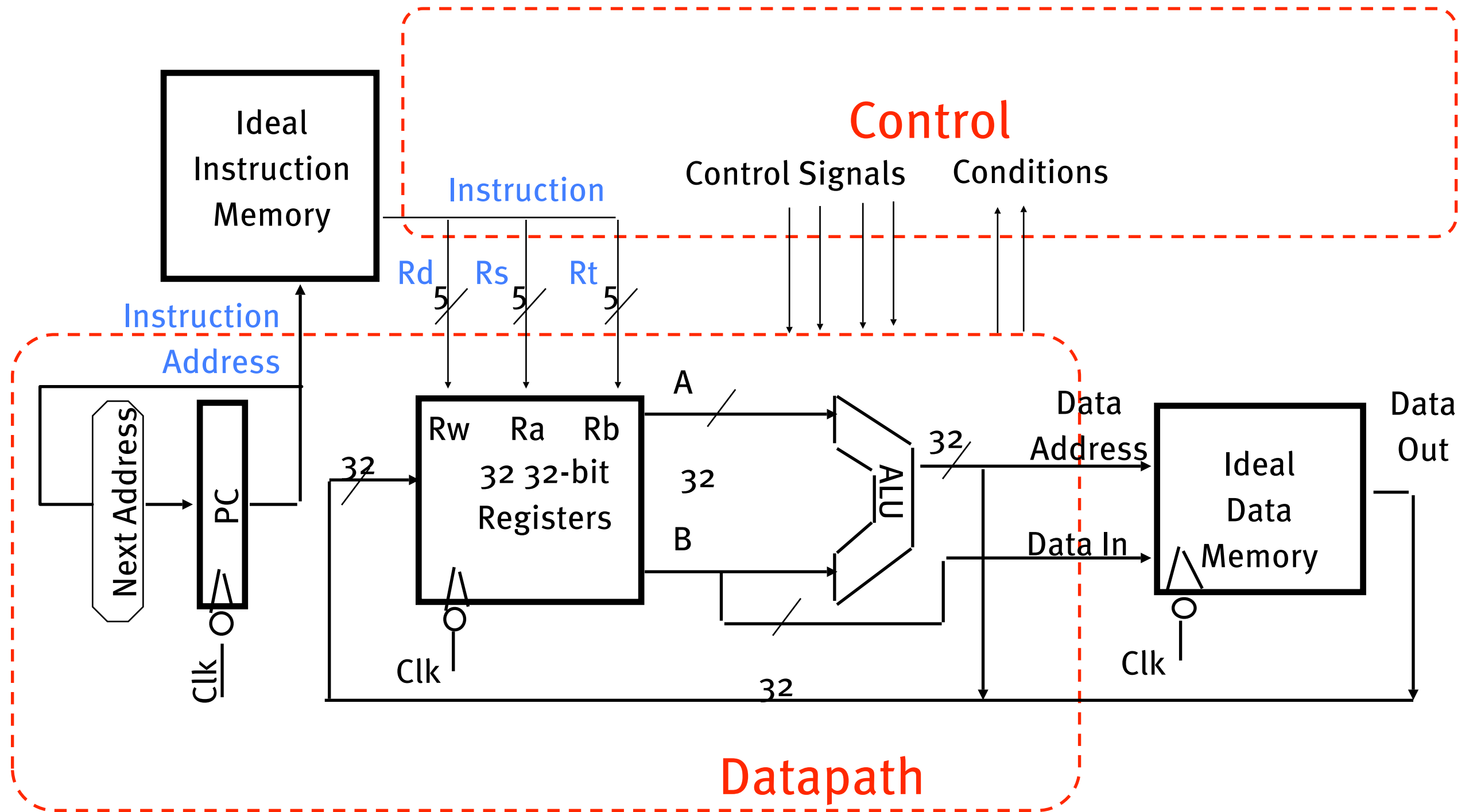


- Old days: ~10 levels of gates
- Today: determined by numerous time-of-flight issues + gate delays
  - clock propagation, wire lengths, drivers

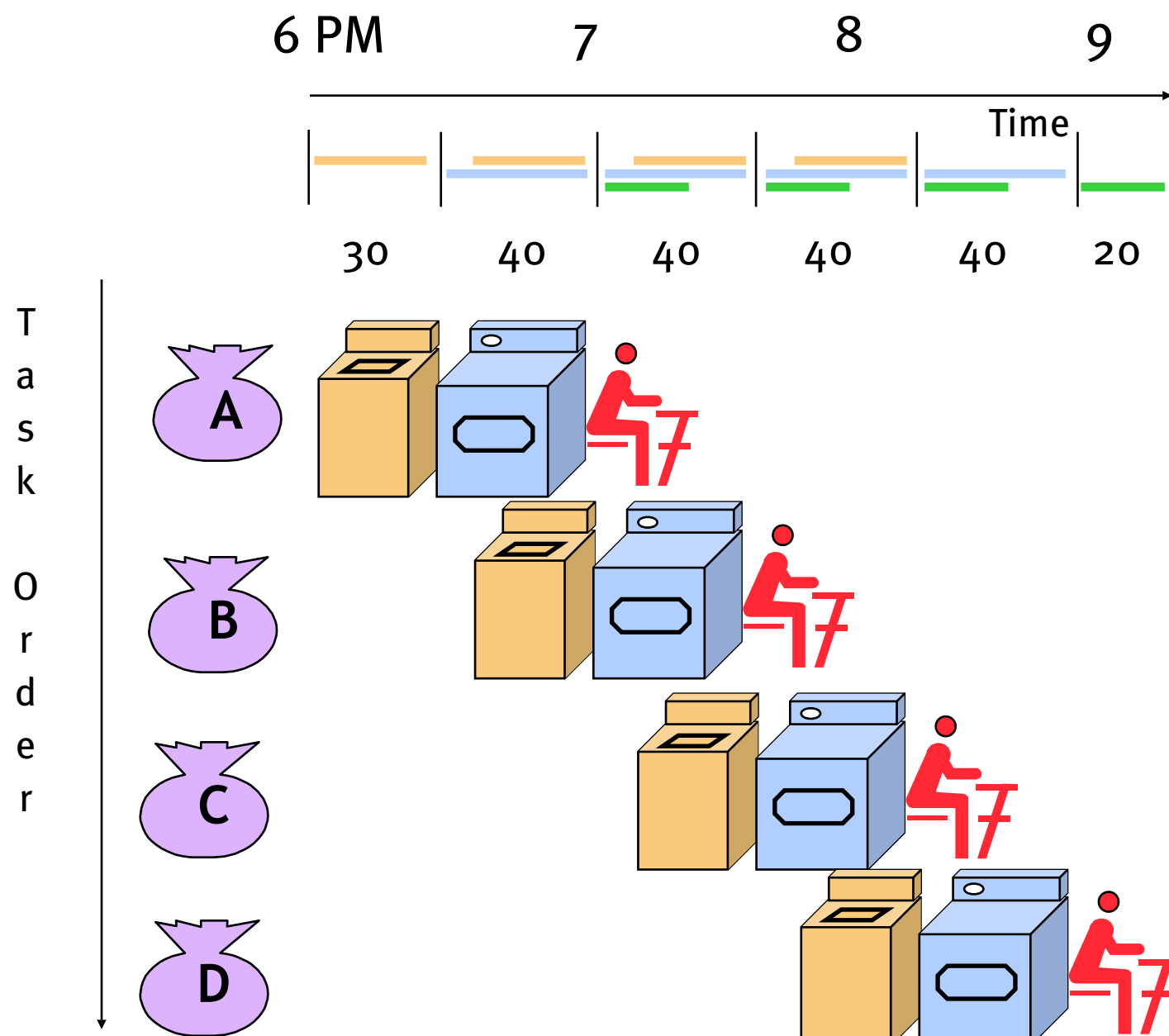
# Putting it All Together: A Single Cycle Datapath



# An Abstract View of Single Cycle

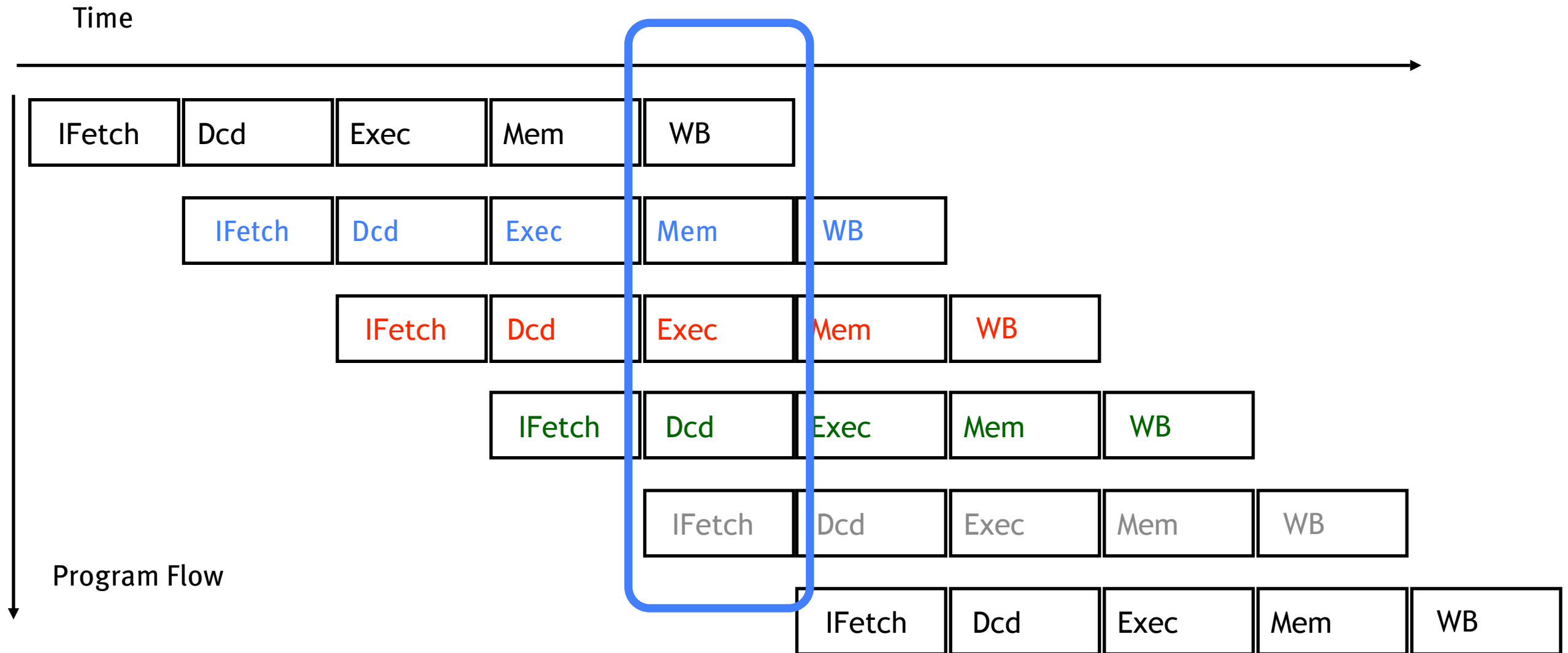


# Pipelining Overview

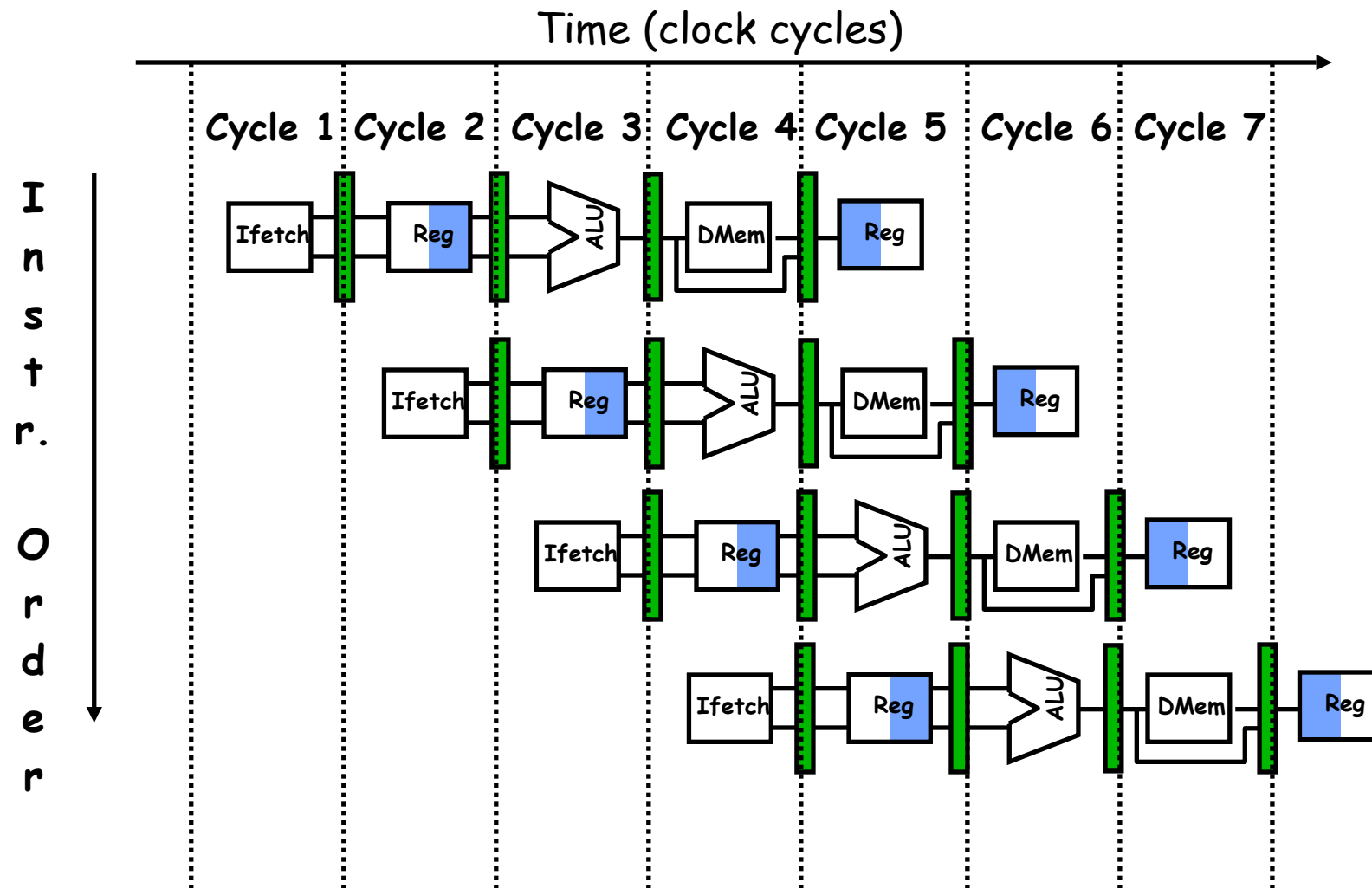


- Pipelining doesn't help latency of single task, it helps throughput of entire workload
- Pipeline rate limited by slowest pipeline stage
- Multiple tasks operating simultaneously using different resources
- Potential speedup = Number pipe stages
- Unbalanced lengths of pipe stages reduces speedup
- Time to "fill" pipeline and time to "drain" it reduces speedup
- Stall for Dependencies

# Conventional Pipelined Execution Representation



# Why Pipeline? Because we can!

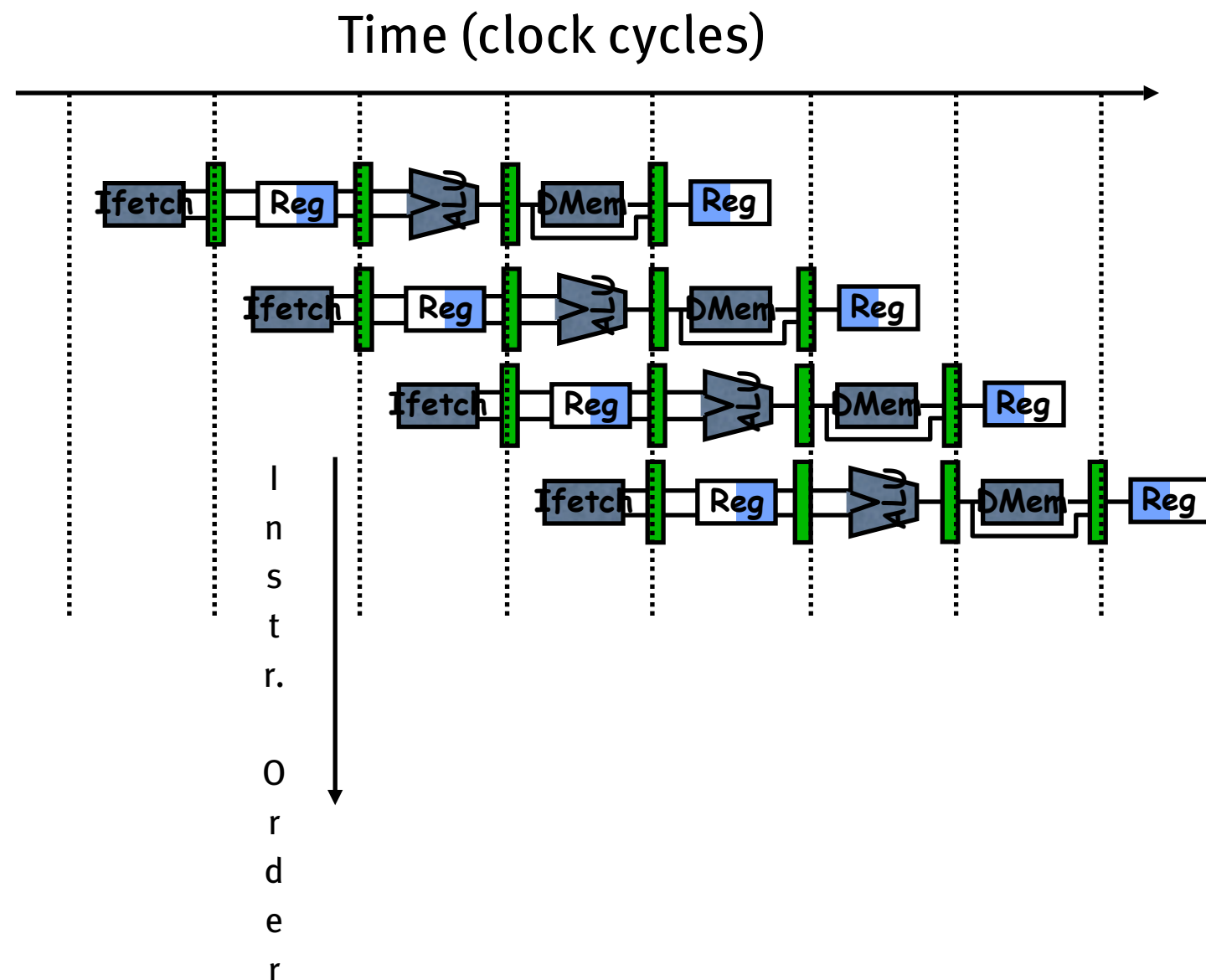


# Why is MIPS great for pipelining?

- All MIPS instructions same length
- Source registers located in same place for every instruction
  - Overlap register fetch and instruction decode
- Simple memory operations
  - MIPS: execute calculates memory address, memory load/store in next stage
  - X86: can operate on result of load: execute calculates memory address, memory load/store in next stage, THEN ALU stage afterwards
- All instructions aligned in memory — 1 access for each instruction

# Limits to pipelining

- **Hazards** prevent next instruction from executing during its designated clock cycle
  - **Structural hazards**: attempt to use the same hardware to do two different things at once
  - **Data hazards**: Instruction depends on result of prior instruction still in the pipeline
  - **Control hazards**: Caused by delay between the fetching of instructions and decisions about changes in control flow (branches and jumps).





# Focus on the Common Case

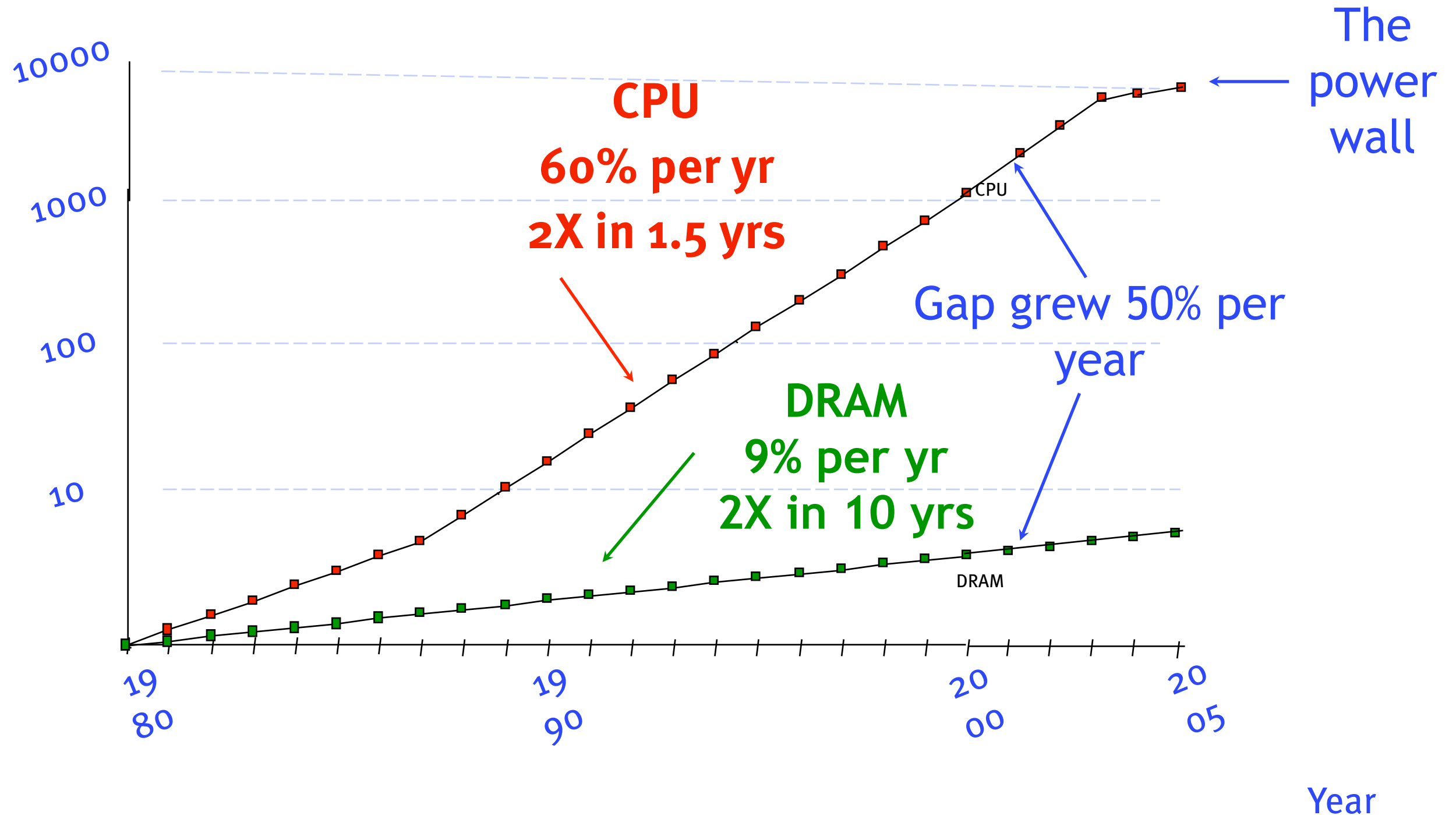
- Common sense guides computer design
  - Since it's engineering, common sense is valuable
- In making a design trade-off, favor the frequent case over the infrequent case
  - e.g., Instruction fetch and decode unit used more frequently than multiplier, so optimize it 1st
  - e.g., If database server has 50 disks / processor, storage dependability dominates system dependability, so optimize it 1st
- Frequent case is often simpler and can be done faster than the infrequent case
  - e.g., overflow is rare when adding 2 numbers, so improve performance by optimizing more for the common case of no overflow
  - May slow down overflow, but overall performance improved by optimizing for the normal case
- What is frequent case and how much performance improved by making case faster?  
⇒ Amdahl's Law

# Pipeline Summary

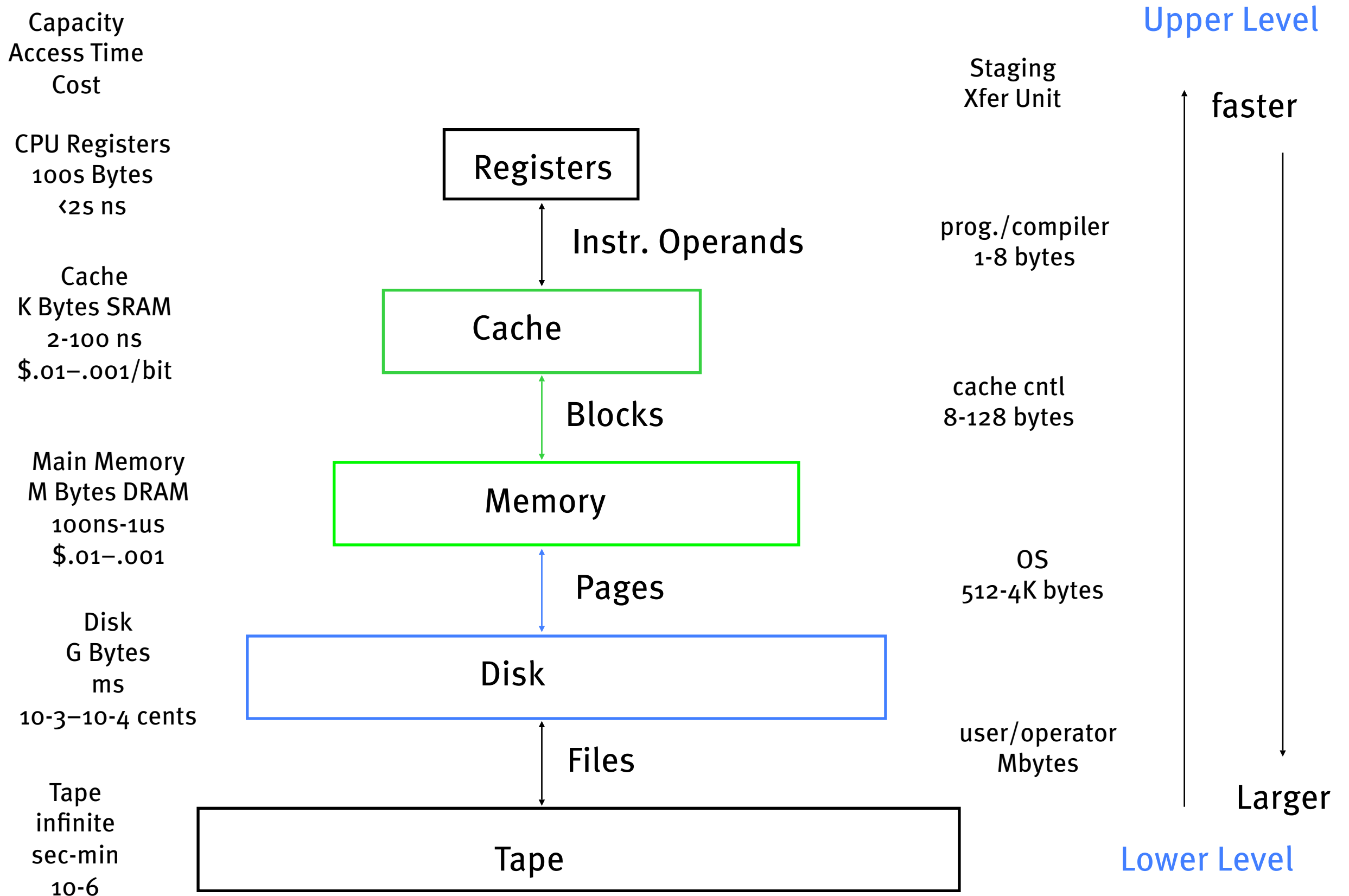
- Simple 5-stage pipeline: F D E M W
- Pipelines pass control information down the pipe just as data moves down pipe
- Resolve data hazards through forwarding.
- Forwarding/Stalls handled by local control
- MIPS I instruction set architecture made pipeline visible (delayed branch, delayed load)
- More performance from deeper pipelines, parallelism

# Why Do We Care About the Memory Hierarchy?

Performance  
(1/latency)



# Levels of the Memory Hierarchy



# Memory Hierarchy

- The Principle of Locality:
  - Program access a relatively small portion of the address space at any instant of time.
    - Temporal Locality: Locality in Time
    - Spatial Locality: Locality in Space
- Three Major Categories of Cache Misses:
  - Compulsory Misses: sad facts of life. Example: cold start misses.
  - Conflict Misses: increase cache size and/or associativity.
  - Capacity Misses: increase cache size

# Design Philosophies Change

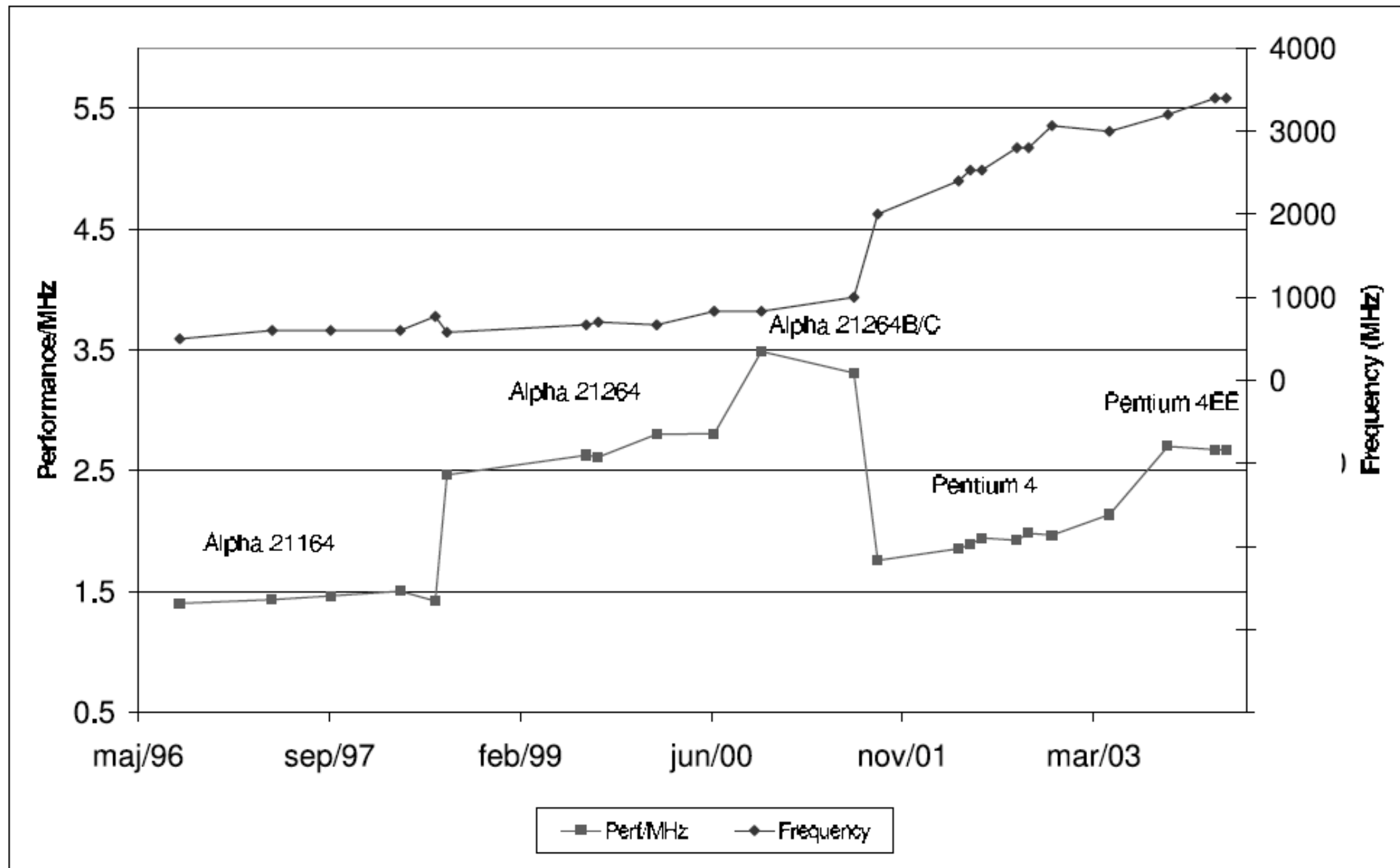


Fig. 3. Performance per MHz 1993-2004.

- Ekman et al., “An In-Depth Look at Computer Performance Growth”

**Looking Forward**

# “Fast enough” ?

- 10 years ago: buy fastest computer you can afford
- Today: No longer the case
- No killer apps!
  - 3D graphics, full-screen video, Internet, speech ...
- You should invent some.



# Gelsinger's Law

- “New generation microarchitectures use twice as many transistors for a 40% increase in performance.”

# Cost of Fabs

- Rock's Law: Cost of fabs double every 4 years
- \$3B for current fab
  - Rise of fabless design houses
  - Rise of for-hire fabs (TSMC, etc.)
  - \$3B fab means \$6B in revenue is required

# Looking To The Future

## The Landscape of Parallel Computing Research: A View from Berkeley



*Krste Asanovic*  
*Ras Bodik*  
*Bryan Christopher Catanzaro*  
*Joseph James Gebis*  
*Parry Husbands*  
*Kurt Keutzer*  
*David A. Patterson*  
*William Lester Plishker*  
*John Shalf*  
*Samuel Webb Williams*  
*Katherine A. Yelick*

Electrical Engineering and Computer Sciences  
University of California at Berkeley

Technical Report No. UCB/EECS-2006-183

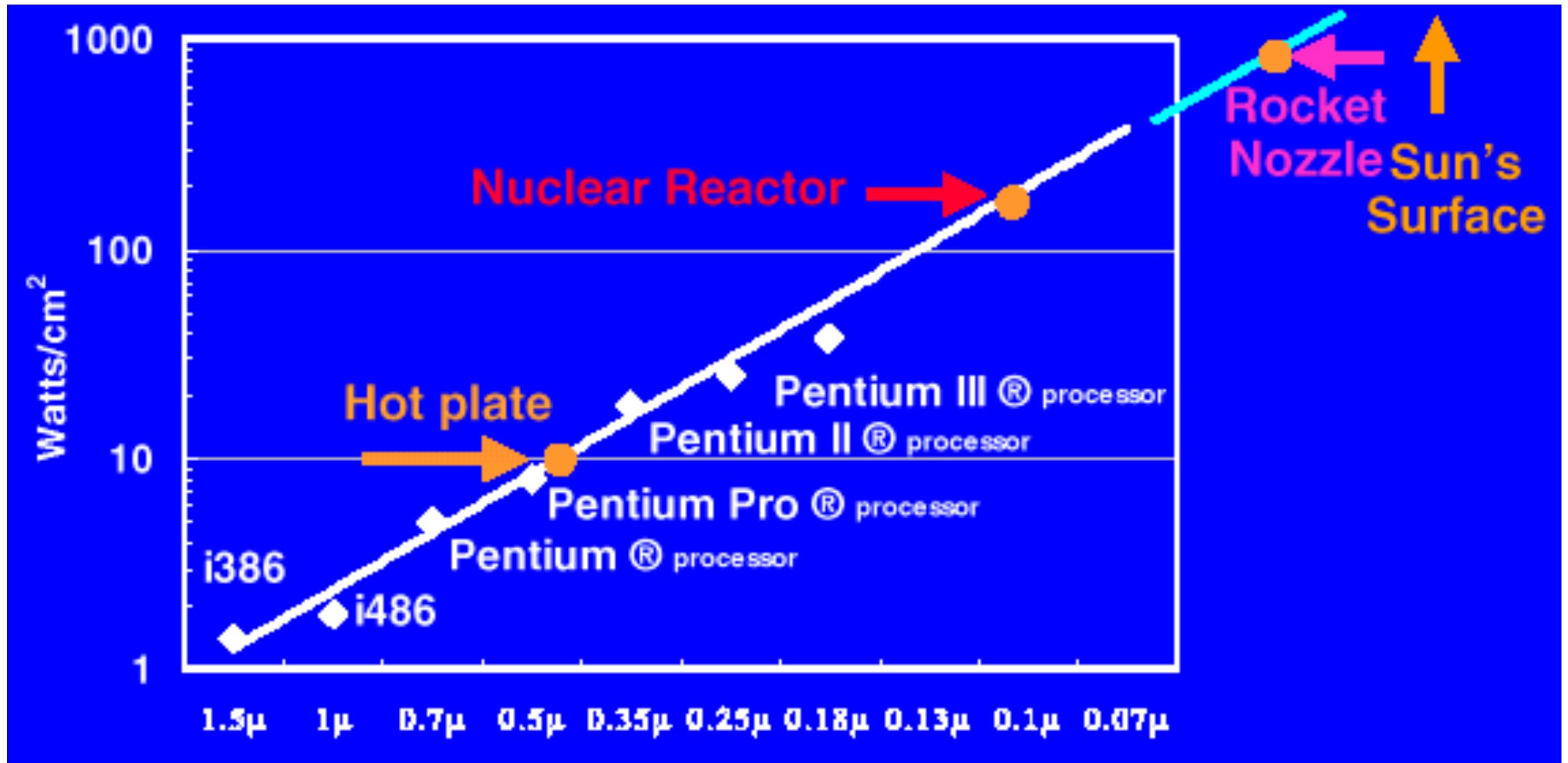
<http://www.eecs.berkeley.edu/Pubs/TechRpts/2006/EECS-2006-183.html>

December 18, 2006

# Power

- Old CW: Power is free, but transistors are expensive.
- New CW is the “Power Wall”: Power is expensive, but transistors are free. We can put more transistors on a chip than we have power to turn on.

# Power Limits Performance



- Bob says: We're pushing perf & clk rates too hard

[courtesy of Bob Colwell]

# Compute vs. Memory

- Old CW: Multiply is slow, but load and store is fast.
- New CW is the *Memory Wall*. Load and store is slow, but multiply is fast. Modern microprocessors can take 200 clocks to access Dynamic Random Access Memory (DRAM), but even floating-point multiplies may take only four clock cycles.

# Computation vs. Communication

- 20 years ago: computation expensive, wires free
  - To first order: ignore wire delay
- Light moves 1 foot/ns in vacuum
  - Wires are also getting thinner
  - Wire delay now significant even on chip!
- Moore's Law implies:
  - Computation gets cheaper
  - Speed of light doesn't change
  - Compute don't communicate!

# Reliability and Test

- Design team size growing at Moore's Law rates
- Processors:
  - Getting larger
  - Getting more complicated
- How to test processors?
- How to prove processor designs correct?



# Intel Directions

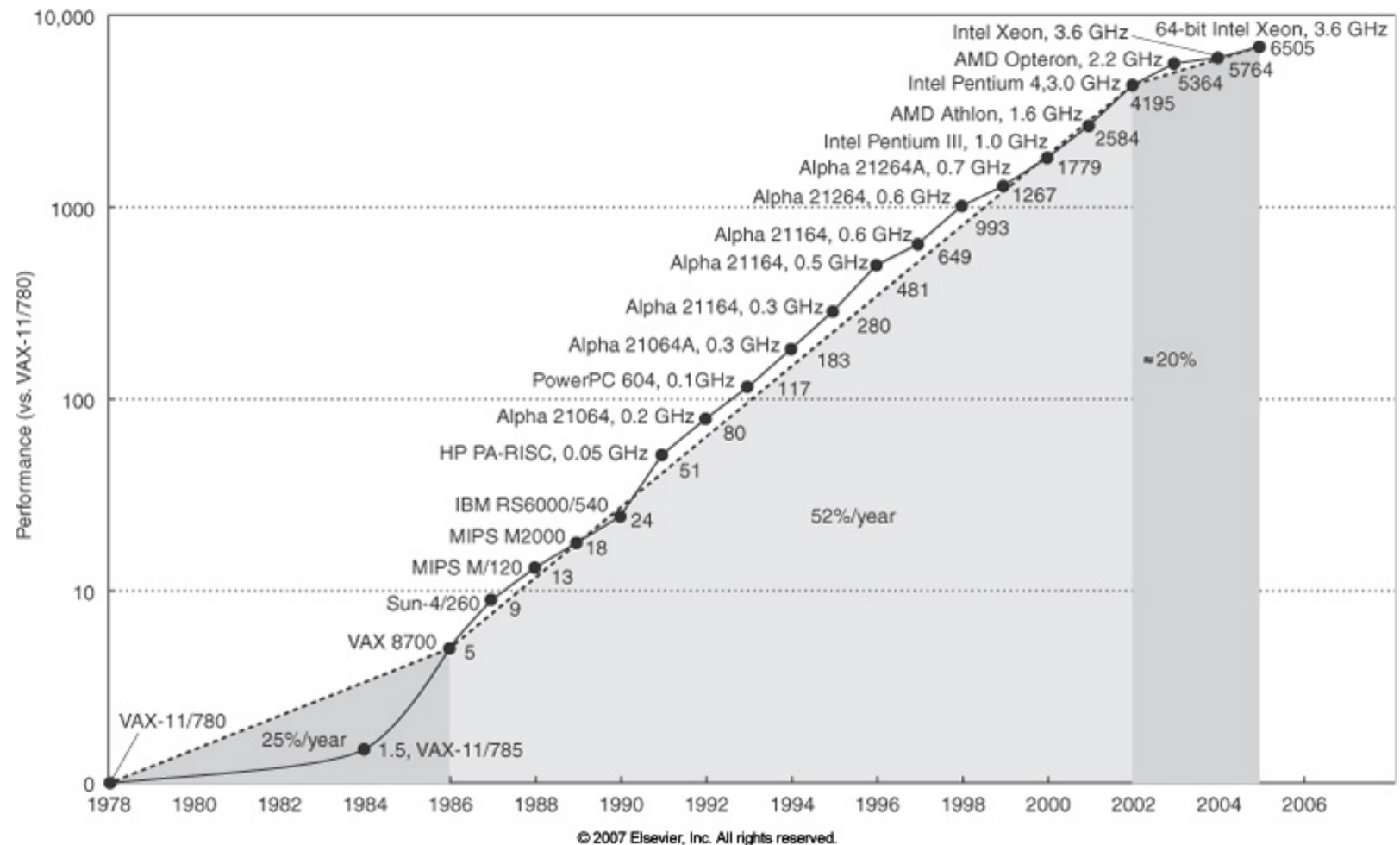
- (Intel Developer Forum, September 2004)
- “What Intel announced at this IDF was no less than a total rethinking of their approach to microprocessors.”
- Moore’s-Law-driven performance scaling will come not from increases not in MHz ratings but in machine width.
  - Power wall
  - Threading
  - MIPS/watt instead of MIPS
- Datasets are growing in size, and so are the network pipes that connect those datasets.
  - Intel claims “doubling of digital data every 18 months”
  - More integration? WiMax?

# Instruction-Level Parallelism

- Old CW: We can reveal more instruction-level parallelism (ILP) via compilers and architecture innovation. Examples from the past include branch prediction, out-of-order execution, speculation, and Very Long Instruction Word systems.
- This is our first three weeks.
- New CW is the “ILP wall”: There are diminishing returns on finding more ILP.

# Uniprocessor Performance

- Old CW: Uniprocessor performance doubles every 18 months.
- New CW is Power Wall + Memory Wall + ILP Wall = Brick Wall. In 2006, performance is a factor of three below the traditional doubling every 18 months that we enjoyed between 1986 and 2002. The doubling of uniprocessor performance may now take 5 years.



# Why EEC 171?

- Old CW: Don't bother parallelizing your application, as you can just wait a little while and run it on a much faster sequential computer.
- New CW: It will be a very long wait for a faster sequential computer.
- Old CW: Increasing clock frequency is the primary method of improving processor performance.
- New CW: Increasing parallelism is the primary method of improving processor performance.
- Old CW: Less than linear scaling for a multiprocessor application is failure.
- New CW: Given the switch to parallel computing, any speedup via parallelism is a success.