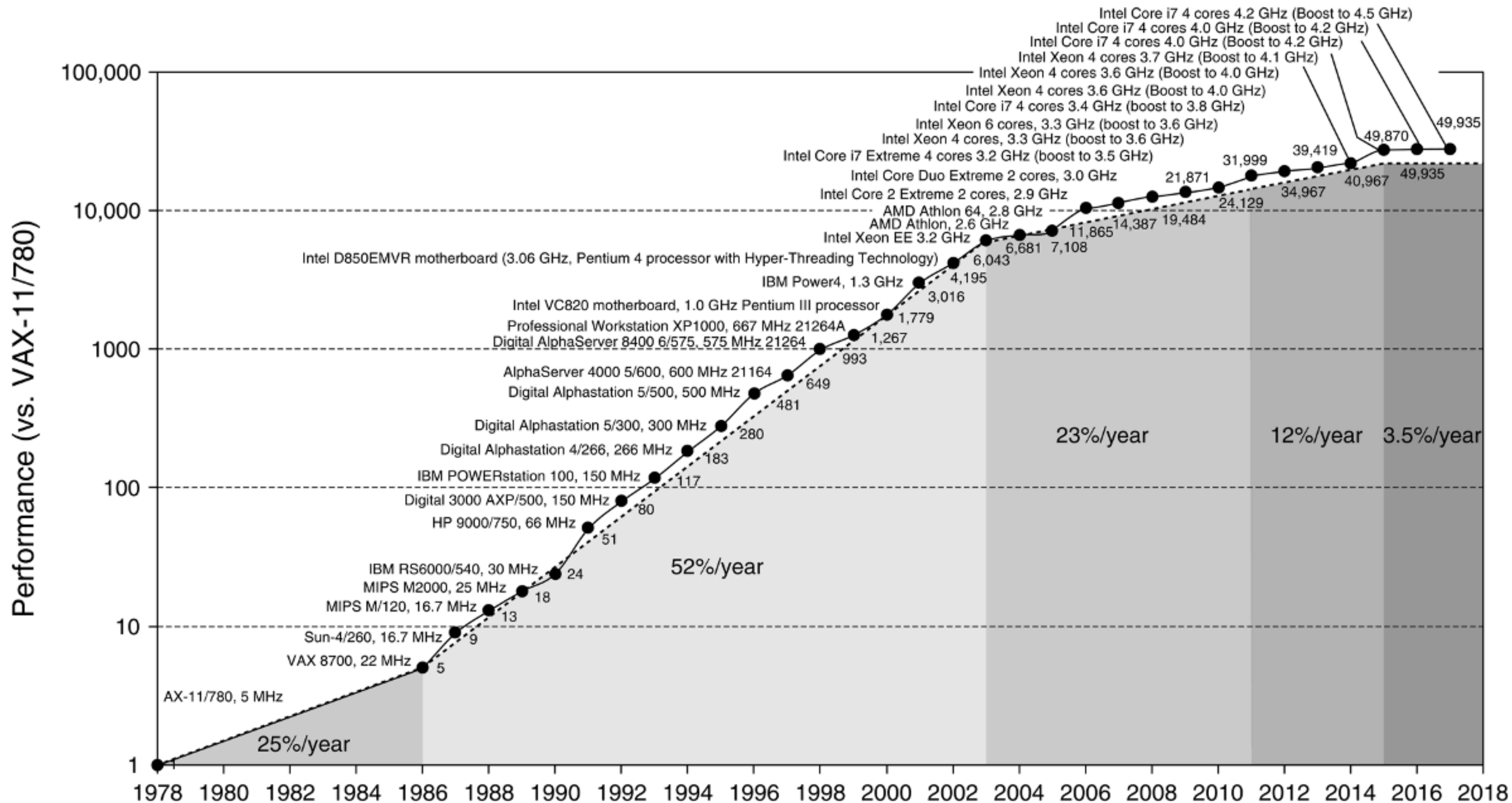


# **Computer Architecture**

## **Chapter 1**



# Classes of Computers

---

- **Personal Mobile Device (PMD)**
  - e.g. smart phones, tablet computers
  - Emphasis on energy efficiency and real-time
- **Desktop Computing**
  - Emphasis on price-performance
- **Servers**
  - Emphasis on availability, scalability, throughput
- **Clusters / Warehouse Scale Computers**
  - Used for “Software as a Service (SaaS)”
  - Emphasis on availability and price-performance
  - Sub-class: Supercomputers, emphasis: floating-point performance and fast internal networks
- **Embedded Computers**
  - Emphasis: price

Feature	Personal mobile device (PMD)	Desktop	Server	Clusters/warehouse-scale computer	Internet of things/embedded
Price of system	\$100–\$1000	\$300–\$2500	\$5000–\$10,000,000	\$100,000–\$200,000,000	\$10–\$100,000
Price of microprocessor	\$10–\$100	\$50–\$500	\$200–\$2000	\$50–\$250	\$0.01–\$100
Critical system design issues	Cost, energy, media performance, responsiveness	Price-performance, energy, graphics performance	Throughput, availability, scalability, energy	Price-performance, throughput, energy proportionality	Price, energy, application-specific performance

**Figure 1.2 A summary of the five mainstream computing classes and their system characteristics.** Sales in 2015 included about 1.6 billion PMDs (90% cell phones), 275 million desktop PCs, and 15 million servers. The total number of embedded processors sold was nearly 19 billion. In total, 14.8 billion ARM-technology-based chips were shipped in 2015. Note the wide range in system price for servers and embedded systems, which go from USB keys to network routers. For servers, this range arises from the need for very large-scale multiprocessor systems for high-end transaction processing.



# Parallelism

---

- **Classes of parallelism in applications:**
  - Data-Level Parallelism (DLP)
  - Task-Level Parallelism (TLP)
- **Classes of architectural parallelism:**
  - Instruction-Level Parallelism (ILP)
  - Vector architectures/Graphic Processor Units (GPUs)
  - Thread-Level Parallelism
  - Request-Level Parallelism



# Flynn's Taxonomy (1966)

---

- Single instruction stream, single data stream (SISD)
- Single instruction stream, multiple data streams (SIMD)
  - Vector architectures
  - Multimedia extensions
  - Graphics processor units
- Multiple instruction streams, single data stream (MISD)
  - No commercial implementation
- Multiple instruction streams, multiple data streams (MIMD)
  - Tightly-coupled MIMD
  - Loosely-coupled MIMD

# Instruction Set Architecture

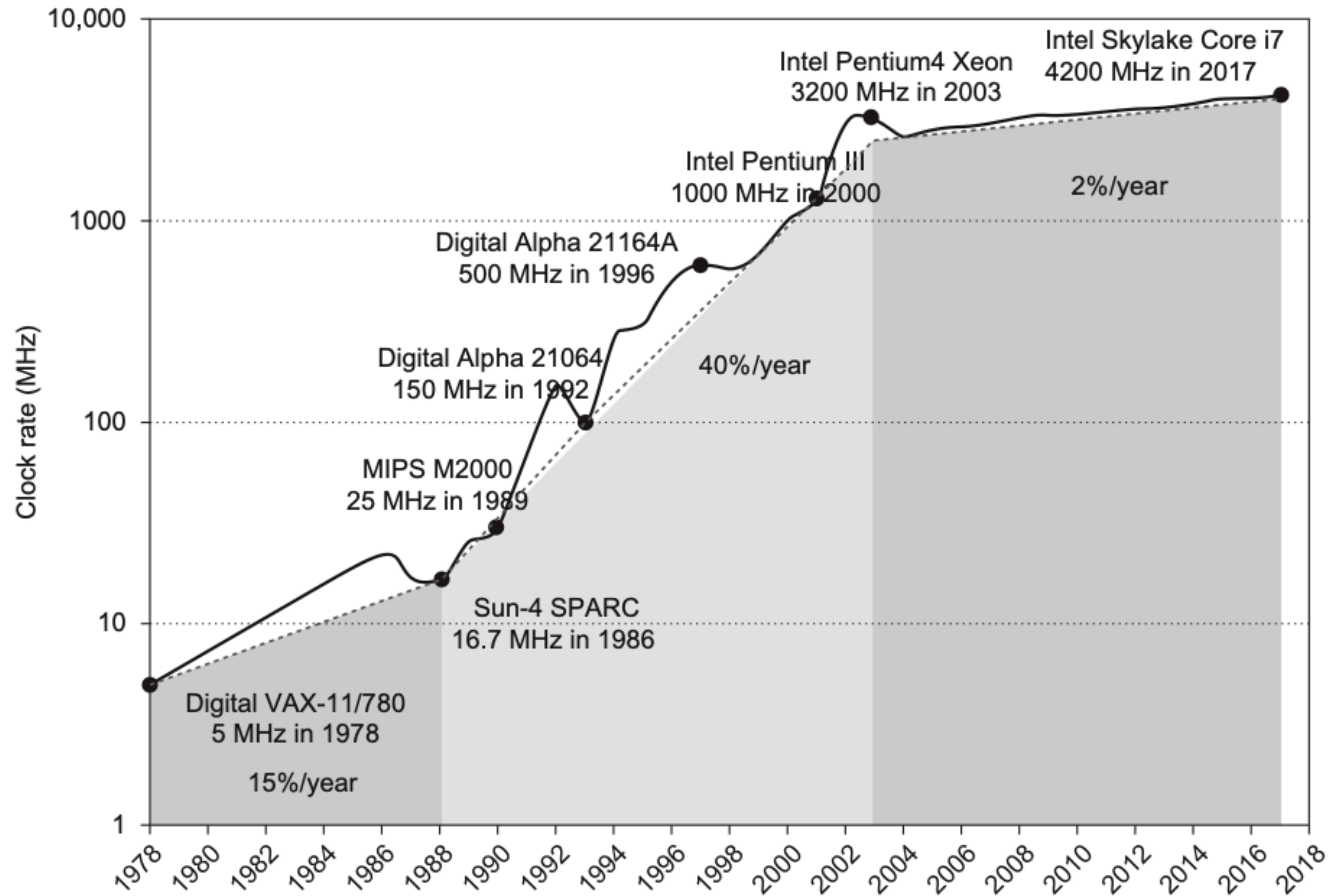
## Some examples

- 80x86
- IBM 360
- VAX 11
- MIPS
- SPARC
- ARM
- RISC-V

# ISA Components

- ISA Classes (Accumulator, Stack, register-memory, load-store)
- Memory addressing (how processor views memory)
- Addressing modes (ways that memory locations can be referenced)
- Operand types, sizes (bit, byte, word, register, integers, floating pt, etc.)
- Operations (arithmetic, logical, compare, branch, jump, system, etc.)
- Control flow (how are compares done, affect on branches, etc.)
- Bit encoding (how are instructions encoded as bit strings)





**Figure 1.11 Growth in clock rate of microprocessors in Figure 1.1.** Between 1978 and 1986, the clock rate improved less than 15% per year while performance improved by 22% per year. During the “renaissance period” of 52% performance improvement per year between 1986 and 2003, clock rates shot up almost 40% per year. Since then, the clock rate has been nearly flat, growing at less than 2% per year, while single processor performance improved recently at just 3.5% per year.



# Measuring Performance

*"Computer X is n times faster than computer Y"*

$$\frac{\text{Execution time}_y}{\text{Execution time}_x} = n$$

$$n = \frac{\text{Execution time}_y}{\text{Execution time}_x} = \frac{\text{Performance}_x}{\text{Performance}_y}$$

CHAP10010





# Benchmarks

- 1. Real Applications*
- 2. Modified (or scripted) applications*
- 3. Kernels*
- 4. Toy benchmarks*
- 5. Synthetic benchmarks*

CHAP1020





	Benchmark name by SPEC generation					
	SPEC2017	SPEC2006	SPEC2000	SPEC95	SPEC92	SPEC89
GNU C compiler	←					gcc
Perl interpreter	←			perl		espresso
Route planning	←		mcf			li
General data compression	XZ		bzip2		compress	eqntott
Discrete Event simulation - computer network	←	omnetpp	vortex	go	sc	
XML to HTML conversion via XSLT	←	xalancbmk	gzip	ijpeg		
Video compression	X264	h264ref	eon	m88ksim		
Artificial Intelligence: alpha-beta tree search (Chess)	deepsjeng	sjeng	twolf			
Artificial Intelligence: Monte Carlo tree search (Go)	leela	gobmk	vortex			
Artificial Intelligence: recursive solution generator (Sudoku)	exchange2	astar	vpr			
		hmmmer	crafty			
		libquantum	parser			
Explosion modeling	←	bwaves				fpppp
Physics: relativity	←	cactuBSSN				tomcatv
Molecular dynamics	←	namd				doduc
Ray tracing	←	povray				nasa7
Fluid dynamics	←	lbm				spice
Weather forecasting	←	wrf			swim	matrix300
Biomedical imaging: optical tomography with finite elements	parest	gamess		apsi	hydro2d	
3D rendering and animation	blender			mgrid	su2cor	
Atmosphere modeling	cam4	milc	wupwise	applu	wave5	
Image manipulation	imagemick	zeusmp	apply	turb3d		
Molecular dynamics	nab	gromacs	galgel			
Computational Electromagnetics	fotonik3d	leslie3d	mesa			
Regional ocean modeling	roms	dealll	art			
		soplex	equake			
		calculix	facerec			
		GemsFDTD	ampp			
		tonto	lucas			
		sphinx3	fma3d			
			sixtrack			



# Quantitative Principles

*1. Make the Common Case Fast*

CHAP1030



University of Idaho





# Quantitative Principles

## *2. Amdahl's Law*

$$\text{Speedup} = \frac{\text{Total Performance using an enhancement}}{\text{Total Performance without the enhancement}}$$

CHAP1040





# Amdahl's Law

$$\text{Speedup} = \frac{\text{Execution time}_{\text{old}}}{\text{Execution time}_{\text{new}}} = \frac{1}{(1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}}}$$

CHAP1065



---

**Example** Suppose that we want to enhance the processor used for web serving. The new processor is 10 times faster on computation in the web serving application than the old processor. Assuming that the original processor is busy with computation 40% of the time and is waiting for I/O 60% of the time, what is the overall speedup gained by incorporating the enhancement?

**Answer**  $\text{Fraction}_{\text{enhanced}} = 0.4; \text{Speedup}_{\text{enhanced}} = 10; \text{Speedup}_{\text{overall}} = \frac{1}{0.6 + \frac{0.4}{10}} = \frac{1}{0.64} \approx 1.56$

---

Amdahl's Law expresses the law of diminishing returns: The incremental improvement in speedup gained by an improvement of just a portion of the computation diminishes as improvements are added. An important corollary of Amdahl's Law is that if an enhancement is usable only for a fraction of a task, then we can't speed up the task by more than the reciprocal of 1 minus that fraction.



# Quantitative Principles

## *3. CPU Performance Equations*

CPU Time = CPU cycles for a program X Clock Cycle Time

or

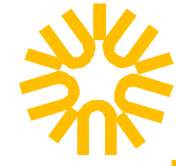
$$\text{CPU time} = \frac{\text{CPU cycles for a program}}{\text{Clock Rate}}$$

$$\text{CPI} = \frac{\text{CPU cycles for a program}}{\text{Instruction Count}}$$

CHAP1050







# Quantitative Principles

## *4. Principle of locality*

*- Temporal Locality*

*- Spacial Locality*

CHAP1060



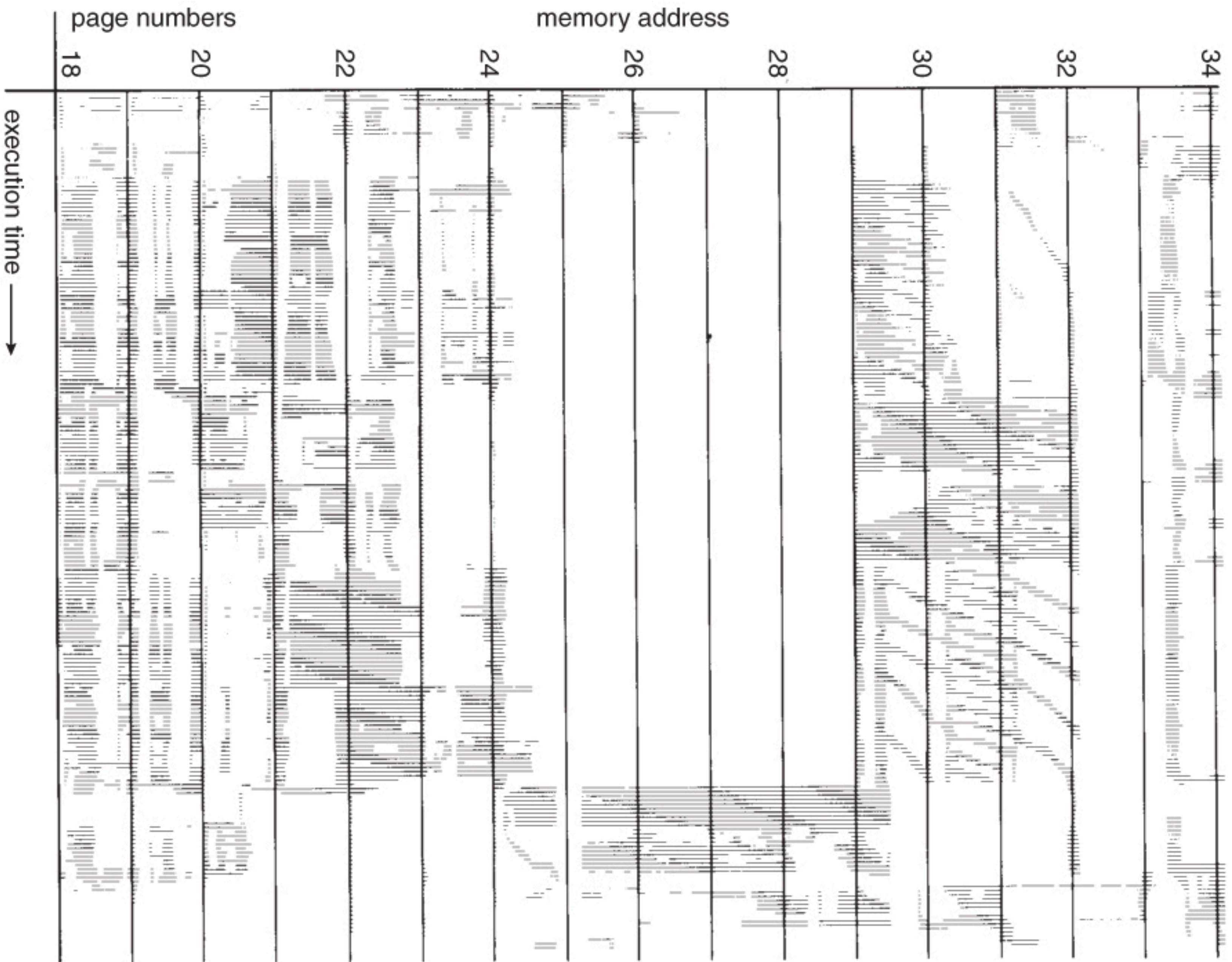


Figure 9.19 - Locality in a memory-reference pattern.



# Quantitative Principles

## *5. Take advantage of parallelism*

CHAP1070



University of Idaho

