

MK COMPUTER ORGANIZATION AND DESIGN
The Hardware/Software Interface **RISC-V Edition**

EECS2021E


Computer Organization Fall 2017

These slides are based on the slides by the authors. The slides doesn't include all the material covered in the lecture. The slides will be explained, modified, and sometime corrected in the lecture.

EECS2021E

- Computer Organization and Design– The hardware/Software approach – RISC-V Edition
- Patterson and Hennessy
- Morgan kaufmann
- Assessment:

■ Participation	5%
■ Quizzes	15%
■ Lab	20%
■ Midterm	25%
■ Final	35%



MK **MK** Chapter 1 — Computer Abstractions and Technology — 2

The Computer Revolution

- Progress in computer technology
 - Underpinned by Moore's Law
- Makes novel applications feasible
 - Computers in automobiles
 - Cell phones
 - Human genome project
 - World Wide Web
 - Search Engines
- Computers are pervasive

MK **MK** Chapter 1 — Computer Abstractions and Technology — 3

Classes of Computers

- Personal computers
 - General purpose, variety of software
 - Subject to cost/performance tradeoff
- Server computers
 - Network based
 - High capacity, performance, reliability
 - Range from small servers to building sized

MK MK Chapter 1 — Computer Abstractions and Technology — 4

Classes of Computers

- Supercomputers
 - High-end scientific and engineering calculations
 - Highest capability but represent a small fraction of the overall computer market
- Embedded computers
 - Hidden as components of systems
 - Stringent power/performance/cost constraints

MK MK Chapter 1 — Computer Abstractions and Technology — 5

The PostPC Era

The graph shows sales trends from 2007 to 2012. The y-axis represents sales volume from 0 to 1400. The x-axis represents years from 2007 to 2012. Four data series are plotted: Tablet (blue line), Smart phone sales (orange line), PC (not including tablet) (green line), and Cell phone (not including smart phone) (red line). Tablet sales peak in 2011 at approximately 1300. Smart phone sales rise sharply from 2010 to 2012, reaching about 700. PC sales remain relatively stable around 300-400. Cell phone sales are the lowest, starting around 100 and rising to about 200 by 2012.

Year	Tablet	Smart phone sales	PC (not including tablet)	Cell phone (not including smart phone)
2007	1000	0	300	100
2008	1100	0	300	100
2009	1000	0	300	100
2010	1300	100	350	150
2011	1300	400	350	150
2012	1100	700	350	200

MK MK Chapter 1 — Computer Abstractions and Technology — 6

The PostPC Era

- Personal Mobile Device (PMD)
 - Battery operated
 - Connects to the Internet
 - Hundreds of dollars
 - Smart phones, tablets, electronic glasses
- Cloud computing
 - Warehouse Scale Computers (WSC)
 - Software as a Service (SaaS)
 - Portion of software run on a PMD and a portion run in the Cloud
 - Amazon and Google

Morgan Kaufmann Publishers logo and footer: Chapter 1 — Computer Abstractions and Technology — 7

What You Will Learn

- How programs are translated into the machine language
 - And how the hardware executes them
- The hardware/software interface
- What determines program performance
 - And how it can be improved
- How hardware designers improve performance

Morgan Kaufmann Publishers logo and footer: Chapter 1 — Computer Abstractions and Technology — 8


Understanding Performance

- Algorithm
 - Determines number of operations executed
- Programming language, compiler, architecture
 - Determine number of machine instructions executed per operation
- Processor and memory system
 - Determine how fast instructions are executed
- I/O system (including OS)
 - Determines how fast I/O operations are executed

Morgan Kaufmann Publishers logo and footer: Chapter 1 — Computer Abstractions and Technology — 9


Eight Great Ideas

- Design for **Moore's Law**
- Use **abstraction** to simplify design
- Make the **common case fast**
- Performance via **parallelism**
- Performance via **pipelining**
- Performance via **prediction**
- **Hierarchy** of memories
- **Dependability** via redundancy



Chapter 1 — Computer Abstractions and Technology — 10

Below Your Program



- Application software
 - Written in high-level language
- System software
 - Compiler: translates HLL code to machine code
 - Operating System: service code
 - Handling input/output
 - Managing memory and storage
 - Scheduling tasks & sharing resources
- Hardware
 - Processor, memory, I/O controllers

Chapter 1 — Computer Abstractions and Technology — 11

Levels of Program Code

- High-level language
 - Level of abstraction closer to problem domain
 - Provides for productivity and portability
- Assembly language
 - Textual representation of instructions
- Hardware representation
 - Binary digits (bits)
 - Encoded instructions and data

```

High-level language program (in C)
swap(int v[], int k)
{
    int temp;
    temp = v[k];
    v[k] = v[k+1];
    v[k+1] = temp;
}

Compiler

Assembly language program (for RISC-V)
swap:  sllw x6, x11, 3
      add x6, x10, x6
      ld  x5, 0(x6)
      ld  x7, 8(x6)
      sd  x7, 0(x6)
      sd  x5, 8(x6)
      jalr x0, 0(x11)

Assembler

Binary machine language program (for RISC-V)
0000000001101011001001100010011
00000000011001010000001100110011
0000000001100110011001100000011
0000000110000110011001110000011
00000000110011001100000110011
00000000010000000100000001100111
            
```

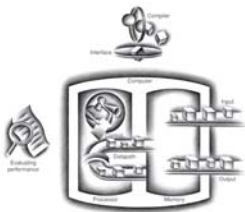
Chapter 1 — Computer Abstractions and Technology — 12

Components of a Computer

3.4 Understanding Computers

The BIG Picture


- Same components for all kinds of computer
 - Desktop, server, embedded
- Input/output includes
 - User-interface devices
 - Display, keyboard, mouse
 - Storage devices
 - Hard disk, CD/DVD, flash
 - Network adapters
 - For communicating with other computers



MK MK Chapter 1 — Computer Abstractions and Technology — 13

Touchscreen

- PostPC device
- Supersedes keyboard and mouse
- Resistive and Capacitive types
 - Most tablets, smart phones use capacitive
 - Capacitive allows multiple touches simultaneously

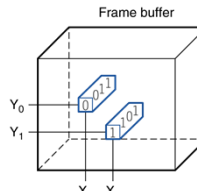


MK MK Chapter 1 — Computer Abstractions and Technology — 14

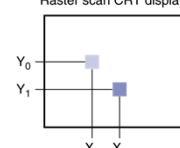
Through the Looking Glass

- LCD screen: picture elements (pixels)
 - Mirrors content of frame buffer memory

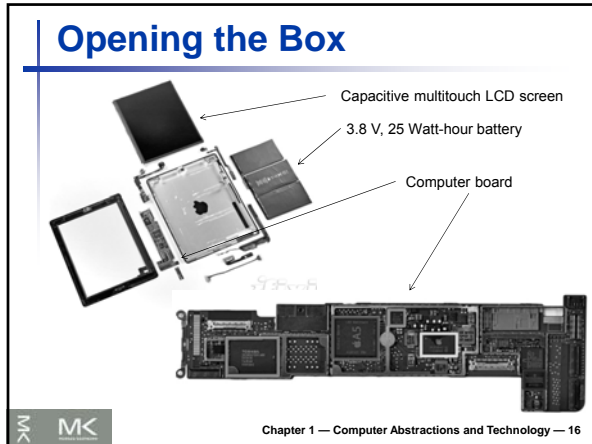
Frame buffer

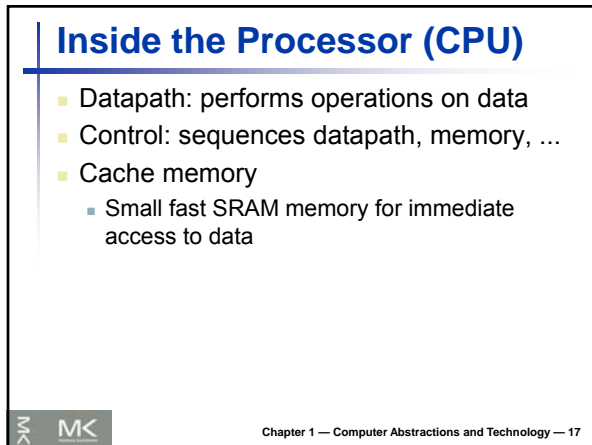


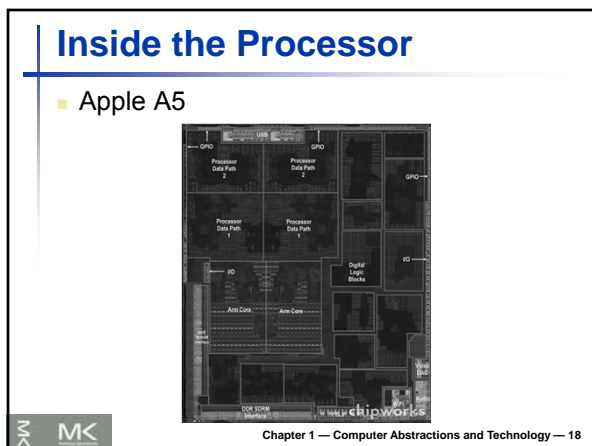
Raster scan CRT display



MK MK Chapter 1 — Computer Abstractions and Technology — 15







Abstractions

The BIG Picture

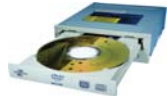
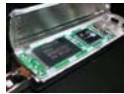
- Abstraction helps us deal with complexity
 - Hide lower-level detail
- Instruction set architecture (ISA)
 - The hardware/software interface
- Application binary interface
 - The ISA plus system software interface
- Implementation
 - The details underlying and interface



Chapter 1 — Computer Abstractions and Technology — 19

A Safe Place for Data

- Volatile main memory
 - Loses instructions and data when power off
- Non-volatile secondary memory
 - Magnetic disk
 - Flash memory
 - Optical disk (CDROM, DVD)



Chapter 1 — Computer Abstractions and Technology — 20

Networks

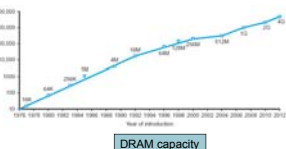
- Communication, resource sharing, nonlocal access
- Local area network (LAN): Ethernet
- Wide area network (WAN): the Internet
- Wireless network: WiFi, Bluetooth



Chapter 1 — Computer Abstractions and Technology — 21

Technology Trends

- Electronics technology continues to evolve
 - Increased capacity and performance
 - Reduced cost



Year	Technology	Relative performance/cost
1951	Vacuum tube	1
1965	Transistor	35
1975	Integrated circuit (IC)	900
1995	Very large scale IC (VLSI)	2,400,000
2013	Ultra large scale IC	250,000,000,000

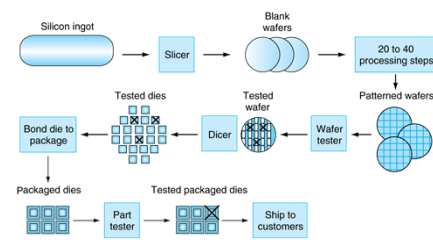
Chapter 1 — Computer Abstractions and Technology — 22

Semiconductor Technology

- Silicon: semiconductor
- Add materials to transform properties:
 - Conductors N or P.
 - Switch Combine them to make switches.

Chapter 1 — Computer Abstractions and Technology — 23

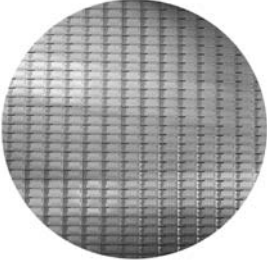
Manufacturing ICs



- Yield: proportion of working dies per wafer

Chapter 1 — Computer Abstractions and Technology — 24

Intel Core i7 Wafer



- 300mm wafer, 280 chips, 32nm technology
- Each chip is 20.7 x 10.5 mm

Morgan Kaufmann Publishers logo and Chapter 1 — Computer Abstractions and Technology — 25

Integrated Circuit Cost

$$\text{Cost per die} = \frac{\text{Cost per wafer}}{\text{Dies per wafer} \times \text{Yield}}$$

$$\text{Dies per wafer} \approx \frac{\text{Wafer area}}{\text{Die area}}$$

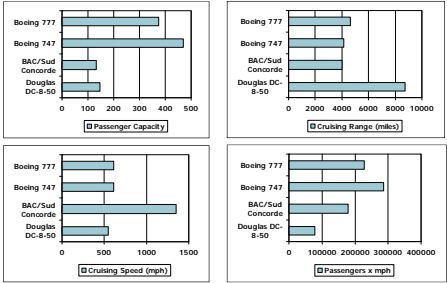
$$\text{Yield} = \frac{1}{(1 + (\text{Defects per area} \times \text{Die area}/2))^2}$$

- Nonlinear relation to area and defect rate
 - Wafer cost and area are fixed
 - Defect rate determined by manufacturing process
 - Die area determined by architecture and circuit design

Morgan Kaufmann Publishers logo and Chapter 1 — Computer Abstractions and Technology — 26

Defining Performance

Which airplane has the best performance?



Morgan Kaufmann Publishers logo and Chapter 1 — Computer Abstractions and Technology — 27

Response Time and Throughput

- Response time
 - How long it takes to do a task
- Throughput
 - Total work done per unit time
 - e.g., tasks/transactions/... per hour
- How are response time and throughput affected by
 - Replacing the processor with a faster version?
 - Adding more processors?
- We'll focus on response time for now...

MK MK Chapter 1 — Computer Abstractions and Technology — 28

Relative Performance

- Define Performance = 1/Execution Time
- "X is n time faster than Y"
 - $Performance_x / Performance_y$
 - $= Execution\ time_y / Execution\ time_x = n$
- Example: time taken to run a program
 - 10s on A, 15s on B
 - $Execution\ Time_B / Execution\ Time_A$
= 15s / 10s = 1.5
 - So A is 1.5 times faster than B

MK MK Chapter 1 — Computer Abstractions and Technology — 29

Measuring Execution Time

- Elapsed time
 - Total response time, including all aspects
 - Processing, I/O, OS overhead, idle time
 - Determines system performance
- CPU time
 - Time spent processing a given job
 - Discounts I/O time, other jobs' shares
 - Comprises user CPU time and system CPU time
 - Different programs are affected differently by CPU and system performance

MK MK Chapter 1 — Computer Abstractions and Technology — 30

CPU Clocking

- Operation of digital hardware governed by a constant-rate clock

- Clock period: duration of a clock cycle
 - e.g., 250ps = 0.25ns = 250×10^{-12} s
- Clock frequency (rate): cycles per second
 - e.g., 4.0GHz = 4000MHz = 4.0×10^9 Hz

Morgan Kaufmann Publishers logo and Chapter 1 — Computer Abstractions and Technology — 31

CPU Time

CPU Time = CPU Clock Cycles \times Clock Cycle Time

$$= \frac{\text{CPU Clock Cycles}}{\text{Clock Rate}}$$

- Performance improved by
 - Reducing number of clock cycles
 - Increasing clock rate
 - Hardware designer must often trade off clock rate against cycle count

Morgan Kaufmann Publishers logo and Chapter 1 — Computer Abstractions and Technology — 32

CPU Time Example

- Computer A: 2GHz clock, 10s CPU time
- Designing Computer B
 - Aim for 6s CPU time
 - Can do faster clock, but causes 1.2 \times clock cycles
- How fast must Computer B clock be?

$$\frac{T_A}{T_B} = \frac{IC_A \times TC_A}{IC_B \times TC_B} = \frac{10}{6}$$

$$\frac{IC_A \times CLOCK_B}{IC_B \times CLOCK_A} = \frac{10}{6}$$

$$\frac{IC_A \times CLOCK_B}{1.2IC_A \times 2G} = \frac{10}{6}$$

$$CLOCK_B = 4GHz$$

Morgan Kaufmann Publishers logo and Chapter 1 — Computer Abstractions and Technology — 33

Instruction Count and CPI

Clock Cycles = Instruction Count × Cycles per Instruction

CPU Time = Instruction Count × CPI × Clock Cycle Time

$$= \frac{\text{Instruction Count} \times \text{CPI}}{\text{Clock Rate}}$$

- Instruction Count for a program
 - Determined by program, ISA and compiler
- Average cycles per instruction
 - Determined by CPU hardware
 - If different instructions have different CPI
 - Average CPI affected by instruction mix

Morgan Kaufmann Publishers logo and Chapter 1 — Computer Abstractions and Technology — 34

CPI Example

- Computer A: Cycle Time = 250ps, CPI = 2.0
- Computer B: Cycle Time = 500ps, CPI = 1.2
- Same ISA
- Which is faster, and by how much?

CPU Time_A = Instruction Count × CPI_A × Cycle Time_A

$$= I \times 2.0 \times 250\text{ps} = I \times 500\text{ps} \leftarrow \text{A is faster...}$$

CPU Time_B = Instruction Count × CPI_B × Cycle Time_B

$$= I \times 1.2 \times 500\text{ps} = I \times 600\text{ps}$$

CPU Time_B = I × 600ps

CPU Time_A = I × 500ps

$$\frac{\text{CPU Time}_B}{\text{CPU Time}_A} = \frac{I \times 600\text{ps}}{I \times 500\text{ps}} = 1.2 \leftarrow \text{...by this much}$$

Morgan Kaufmann Publishers logo and Chapter 1 — Computer Abstractions and Technology — 35

CPI in More Detail

- If different instruction classes take different numbers of cycles

$$\text{Clock Cycles} = \sum_{i=1}^n (\text{CPI}_i \times \text{Instruction Count}_i)$$

- Weighted average CPI

$$\text{CPI} = \frac{\text{Clock Cycles}}{\text{Instruction Count}} = \sum_{i=1}^n \left(\underbrace{\text{CPI}_i \times \frac{\text{Instruction Count}_i}{\text{Instruction Count}}}_{\text{Relative frequency}} \right)$$

Morgan Kaufmann Publishers logo and Chapter 1 — Computer Abstractions and Technology — 36

CPI Example

- Alternative compiled code sequences using instructions in classes A, B, C

Class	A	B	C
CPI for class	1	2	3
IC in sequence 1	2	1	2
IC in sequence 2	4	1	1

- Sequence 1: IC = 5
 - Clock Cycles = $2 \times 1 + 1 \times 2 + 2 \times 3 = 10$
 - Avg. CPI = $10/5 = 2.0$
- Sequence 2: IC = 6
 - Clock Cycles = $4 \times 1 + 1 \times 2 + 1 \times 3 = 9$
 - Avg. CPI = $9/6 = 1.5$

Chapter 1 — Computer Abstractions and Technology — 37

Performance Summary


The BIG Picture

$$\text{CPU Time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock cycle}}$$

- Performance depends on
 - Algorithm: affects IC, possibly CPI
 - Programming language: affects IC, CPI
 - Compiler: affects IC, CPI
 - Instruction set architecture: affects IC, CPI, T_c

Chapter 1 — Computer Abstractions and Technology — 38

Power Trends



- In CMOS IC technology

$$\text{Power} = \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency}$$

Annotations: $\times 30$, $5V \rightarrow 1V$, $\times 1000$

Chapter 1 — Computer Abstractions and Technology — 39

Reducing Power

- Suppose a new CPU has
 - 85% of capacitive load of old CPU
 - 15% voltage and 15% frequency reduction

$$\frac{P_{\text{new}}}{P_{\text{old}}} = \frac{C_{\text{old}} \times 0.85 \times (V_{\text{old}} \times 0.85)^2 \times F_{\text{old}} \times 0.85}{C_{\text{old}} \times V_{\text{old}}^2 \times F_{\text{old}}} = 0.85^4 = 0.52$$

- The power wall
 - We can't reduce voltage further
 - We can't remove more heat
- How else can we improve performance?

Chapter 1 — Computer Abstractions and Technology — 40

Uniprocessor Performance

Chapter 1 — Computer Abstractions and Technology — 41

Multiprocessors

- Multicore microprocessors
 - More than one processor per chip
- Requires explicitly parallel programming
 - Compare with instruction level parallelism
 - Hardware executes multiple instructions at once
 - Hidden from the programmer
 - Hard to do
 - Programming for performance
 - Load balancing
 - Optimizing communication and synchronization

Chapter 1 — Computer Abstractions and Technology — 42

SPEC CPU Benchmark

- Programs used to measure performance
 - Supposedly typical of actual workload
- Standard Performance Evaluation Corp (SPEC)
 - Develops benchmarks for CPU, I/O, Web, ...
- SPEC CPU2006
 - Elapsed time to execute a selection of programs
 - Negligible I/O, so focuses on CPU performance
 - Normalize relative to reference machine
 - Summarize as geometric mean of performance ratios
 - CINT2006 (integer) and CFP2006 (floating-point)

$$\sqrt[n]{\prod_{i=1}^n \text{Execution time ratio}_i}$$

Chapter 1 — Computer Abstractions and Technology — 43

CINT2006 for Intel Core i7 920

Description	Name	Instruction Count x 10 ³	CPI	Clock cycle time (seconds x 10 ⁻⁹)	Execution Time (seconds)	Reference Time (seconds)	SPECratio
Interpreted string processing	perl	2252	0.60	0.376	908	9770	19.2
Block sorting	tbody2	2390	0.70	0.376	629	9550	15.4
compression							
GNU C compiler	gcc	794	1.20	0.376	358	8050	22.5
Combinatorial optimization	mcf	221	2.66	0.376	221	9120	41.2
Go game (k)	go	1274	1.10	0.376	527	10490	19.9
Search gene sequence	primr	2616	0.60	0.376	590	9330	15.8
Chess game (k)	sjeng	1946	0.80	0.376	565	12100	20.7
Quantum computer simulation	libquantum	659	0.44	0.376	109	20720	190.0
Video compression	h264enc	3793	0.50	0.376	713	22130	31.0
Discrete event simulation library	onnwtp	367	2.10	0.376	290	6250	21.5
Games/path finding	sstar	1250	1.00	0.376	470	7020	14.9
XML parsing	xalanbmh	1045	0.70	0.376	275	6900	25.1
Geometric mean	-	-	-	-	-	-	25.7

Chapter 1 — Computer Abstractions and Technology — 44

SPEC Power Benchmark

- Power consumption of server at different workload levels
 - Performance: ssj_ops/sec
 - Power: Watts (Joules/sec)

$$\text{Overall ssj_ops per Watt} = \frac{\left(\sum_{i=0}^{10} \text{ssj_ops}_i\right)}{\left(\sum_{i=0}^{10} \text{power}_i\right)}$$

Chapter 1 — Computer Abstractions and Technology — 45

SPECpower_ssjs2008 for Xeon X5650

Target Load %	Performance (ssj_ops)	Average Power (Watts)
100%	865,618	258
90%	786,688	242
80%	698,051	224
70%	607,826	204
60%	521,391	185
50%	436,757	170
40%	345,919	157
30%	262,071	146
20%	176,051	135
10%	86,784	121
0%	0	80
Overall Sum	4,787,166	1,922
$\Sigma \text{ssj_ops} / \Sigma \text{power} =$		2,490

Chapter 1 — Computer Abstractions and Technology — 46

Pitfall: Amdahl's Law

- Improving an aspect of a computer and expecting a proportional improvement in overall performance

$$T_{\text{improved}} = \frac{T_{\text{affected}}}{\text{improvement factor}} + T_{\text{unaffected}}$$

- Example: multiply accounts for 80s/100s
 - How much improvement in multiply performance to get 5x overall?

$$20 = \frac{80}{n} + 20 \quad \text{Can't be done!}$$
- Corollary: make the common case fast

Chapter 1 — Computer Abstractions and Technology — 47

Fallacy: Low Power at Idle

- Look back at i7 power benchmark
 - At 100% load: 258W
 - At 50% load: 170W (66%)
 - At 10% load: 121W (47%)
- Google data center
 - Mostly operates at 10% – 50% load
 - At 100% load less than 1% of the time
- Consider designing processors to make power proportional to load

Chapter 1 — Computer Abstractions and Technology — 48

Pitfall: MIPS as a Performance Metric

- MIPS: Millions of Instructions Per Second
 - Doesn't account for
 - Differences in ISAs between computers
 - Differences in complexity between instructions

$$\begin{aligned} \text{MIPS} &= \frac{\text{Instruction count}}{\text{Execution time} \times 10^6} \\ &= \frac{\text{Instruction count}}{\frac{\text{Instruction count} \times \text{CPI}}{\text{Clock rate}} \times 10^6} = \frac{\text{Clock rate}}{\text{CPI} \times 10^6} \end{aligned}$$

- CPI varies between programs on a given CPU



Concluding Remarks

- Cost/performance is improving
 - Due to underlying technology development
- Hierarchical layers of abstraction
 - In both hardware and software
- Instruction set architecture
 - The hardware/software interface
- Execution time: the best performance measure
- Power is a limiting factor
 - Use parallelism to improve performance