Fiscal Year 2024

Ver. 2024-12-09a

Course number: CSC.T433 School of Computing, Graduate major in Computer Science

# Advanced Computer Architecture

1. Design and Analysis of Computer Systems

www.arch.cs.titech.ac.jp/lecture/ACA/ Room No. W8E-308, Lecture (Face-to-face) Mon 13:30-15:10, Thr 13:30-15:10

Kenji Kise, Department of Computer Science kise \_at\_ c.titech.ac.jp

## Syllabus (1/3)

#### Course description and aims

This course aims to provide students with <u>cutting-edge technologies and future trends of computer architecture with focusing on a</u> <u>microprocessor which plays an important role</u> in the downsizing, personalization, and improvement of performance and power consumption of computer systems such as PCs, personal mobile devices, and embedded systems.

In this course, first, along with important concepts of computer architecture, students will learn from instruction set architectures to mechanisms for extracting instruction level parallelism used in out-of-order superscalar processors. After that, students will learn mechanisms for exploiting thread level parallelism adopted in multi-processors and multi-core processors.

#### Student learning outcomes

By taking this course, students will learn:

- (1) Basic principles for building today's high-performance computer systems
- (2) Mechanisms for extracting instruction level parallelism used in high-performance microprocessors
- (3) Methods for exploiting thread level parallelism adopted in multi-processors and multi-core processors
- (4) New inter-relationship between software and hardware

#### Keywords

Computer Architecture, Processor, Embedded System, multi-processor, multi-core processor

#### Competencies that will be developed

Intercultural skills

Communication skills Critic

Critical thinking skills

Practical and/or problem-solving skills

#### **Class flow**

Before coming to class, students should read the course schedule and check what topics will be covered. Required learning should be completed outside of the classroom for preparation and review purposes.



# Syllabus (2/3)

#### Textbook(s)

John L. Hennessy, David A. Patterson. Computer Architecture A Quantitative Approach, Fifth Edition. Morgan Kaufmann Publishers Inc., 2012

#### Reference books, course materials, etc.

William James Dally, Brian Patrick Towles. Principles and Practices of Interconnection Networks. Morgan Kaufman Publishers Inc., 2004.

#### Assessment criteria and methods

Students will be assessed on their understanding of instruction level parallelism, multi-processor, and thread level parallelism. Students' course scores are based on the mid-term report and assignments (40%), and the final report (60%).

#### **Related courses**

CSC.T363 : Computer Architecture

CSC.T341 : Computer Logic Design

Prerequisites (i.e., required knowledge, skills, courses, etc.)

No prerequisites are necessary, but enrollment in the related courses is desirable.

Contact information (e-mail and phone) Notice : Please replace from "[at]" to "@"(half-width character).

Kise Kenji: kise[at]c.titech.ac.jp

#### Office hours

Contact by e-mail in advance to schedule an appointment.



# Support page https://www.arch.cs.titech.ac.jp/lecture/ACA/

SC.T433 Advanced Computer Architecture Support Page, Dept. of Computer S	Science, Science Tokyo
ews	
• This page is updated. (2024-12-09)	
• This page is open. (2018-11-28)	
ecture Slides and Materials	
• Lecture01 2024-12-09 13:30-15:10: Introduction, Design and Analysis of Computer Systems	
Lecture02 2024-12-12 13:30-15:10: Instruction Set Architecture	
• Lecture03 2024-12-16 13:30-15:10: HDL, single-cycle processor, Memory Hierarchy Design	
• Lecture04 2024-12-20 13:30-15:10: Pipelining	
• Lecture05 2024-12-23 13:30-15:10: Instruction Level Parallelism: Concepts and Challenges	
• Lecture06 2024-12-26 13:30-15:10: Instruction Level Parallelism: Instruction Fetch and Brand	ch Prediction
Year-end holidays	
Lecture06 2025-01-06 13:30-15:10: Instruction Level Parallelism: Dynamic Scheduling	
• Lecture07 2025-01-09 13:30-15:10: Instruction Level Parallelism: Exploiting ILP Using Multipl	le Issue and Speculation
• Lecture08 2025-01-16 13:30-15:10: Instruction Level Parallelism: Out-of-order Execution and	l Multithreading
• Lecture09 2025-01-20 13:30-15:10: Multi-Processor: Distributed Memory and Shared Memory	y Architecture
Lecture10 2025-01-23 13:30-15:10: Thread Level Parallelism: Interconnection Network	
Lecture 11 2025-01-27 13:30-15:10: Thread Level Parallelism: Coherence and Synchronization	n
Lecture12 2025-01-30 13:30-15:10: Thread Level Parallelism: Memory Consistency Model	
• Lecture13 2025-02-03	
Final Report by February 10, 2025	
or Verilog HDL description, please refer to <u>lectures 2 to 4 of Computer Logic Design</u> .	



## The birth of microprocessors in 1971





### Growth in processor performance



### Program in C and RISC-V assembly code

```
int main(){
    int i, sum=1;
    for(i=1; i<10; i++) sum = sum + i;
    return sum;
}</pre>
```

main1.c

main1.s

main:			
addi	a0, zero, 1	# sum = 1	
addi	a5, zero, 1	# i = 1	
addi	a4, zero, 10	# a4 = 10	
.L2:			
beq	a5, a4, .L3	# branch if i==10 to .L3	
add	a0, a0, a5	# sum = sum + i	
addi	a5, a5, 1	# i++	
beq	zero, zero, .L2	# branch .L2	
.L3:			
ret		# return	



### Venus RISC-V editor and simulator

#### https://venus.cs61c.org/

♥ ♥ venus	× +					- 0 ×
$\leftrightarrow$ $\rightarrow$ C $\widehat{m}$	venus.cs61c.org				☆ 🗟 🖸   🔇	更新を完了 :
			Venus Editor Simulator Chocopy			
	Run	Prev Reset Dump	Trace Re-assemble from Editor		Registers Memory Cache VDB	<b>A</b>
	- Null Step	They heset Dump			Integer (R) Floating (F)	
PC	Machine Code	Basic Code	Original Code	zero	0×00000000	
0×0	0x00100513	addi x10 x0 1	addi a0, zero, 1 # sum = 1	ra (x1)	0×00000000	
0x4	0x00100793	addi x15 x0 1	addi a5, zero, 1 # i = 1	sp (x2)	0x7FFFFDC	
0×8	0x00A00713	addi x14 x0 10	addi a4, zero, 10 # a4 = 10	gp (x3)	0x10000000	
Øxc	0x00E78863	beq x15 x14 16	beq a5, a4, .L3 # branch if i==10 to .L3	tp (x4)	0×00000000	
0×10	0×00F50533	add x10 x10 x15	add a0, a0, a5 # sum = sum + i	t0 (x5)	0×00000000	
0x14	0x00178793	addi x15 x15 1	addi a5, a5, 1 # i++	t1 (x6)	0×00000000	
0×18	0×FE000AE3	beq x0 x0 -12	beq zero, zero, .L2 # branch .L2	t2 (x7)	0×00000000	
Øx1c	0×00008067	jalr x0 x1 0	ret # return	s0 (x8)	0×00000000	
				s1 (x9)	0x0000000	
				a0 (x10)	0x0000001	
				a1 (x11)	0x7FFFFDC	
				a2 (x12)	0×00000000	
				a3 (x13)	0×00000000	
				a4 (x14)	0×00000000	
				a5 (x15)	0×00000000	
console output	:	Copy! Downloa	d! Clear!	a6 (x16)	0x00000000	
				a7 (x17)	0x00000000	
				Display Settings	Hex 🗸	
				li		



## Defining performance and speedup

- Normally interested in reducing
  - Execution time (response time) the time between the start and the completion of a program
  - Performance is the inverse of execution time.

 $performance_{x} = 1 / execution_time_{x}$ 

• Thus, to maximize performance, need to minimize execution time If X is n times faster than Y, then

\_\_\_\_\_performance<sub>X</sub> = \_\_\_\_\_execution\_time<sub>y</sub> \_\_\_\_\_performance<sub>y</sub> = \_\_\_\_\_execution\_time<sub>X</sub> = n



## Performance Factors

- Want to distinguish elapsed time and the time spent on our task
- CPU execution time (CPU time) : time the CPU spends working on a task
  - Does not include time waiting for I/O or running other programs

CPU execution time	= # CPU clock cycles
for a program	for a program x clock cycle time
	or
CPU execution time	= # CPU clock cycles for a program
for a program	clock rate
CPU execution time	= # CPU clock cycles for a program
for a program	F
Clock rate or Freque	icy F is the inverse of clock cycle time



Multiply both sides if the equation by N Performance\_x\_N = F  $\times N / \# CPU$  clock cycles for a program

- Performance = F x IPC
  - F: frequency (clock rate)
  - IPC : executed instructions per cycle (N / clock cycles)
- The performance can be improved by increasing either F or IPC

The past, present, and future of the world's most important device



#### The Device That Changed Everything

The best

contact

explanation

of the point-

transistor is

in Bardeen's

1956 Nobel

Prize lecture.

that left out

important

details.

but even

Transistors are civilization's invisible infrastructure

EDITOR'S NOTE

was roaming around the IEEE Spectrum office a couple of months ago, looking at the display cases the IEEE History Center has installed in the corridor that runs along the conference rooms at 3 Park. They feature photos of illustrious engineers, plaques for IEEE milestones, and a handful of vintage electronics and memorrabilita including an original Sony Walkman, an Edison Mazda lightbulb, and an RCA Radiotron vacuum tube. And, to my utter surprise and delight, a replica of the first point-contact transistor invented by John Bardeen, Walter Brattain, and William Shockley 75 years ago this month.

I dashed over to our photography director, Randi Klett, and startled her with my excitement, which, when she saw my discovery, she understood: We needed a picture of that replica, which she expertly shot and now accompanies this column.

What amazed me most besides the fact that the very thing this issue is devoted to was here with us? I'd passed by it countless times and never noticed it, even though it is tens of billions times the size of one of today's transistors. In fact, each of us is surrounded by billions, if not trillions of transistors, none of which are visible to the naked eye. It is a testament to imagination and ingenuity of three generations of electronics engineers who took the (by today's standards) mammoth point-contact transistors and shrunk it down to the point where transistors are so ubiquitous that civilization as we know it would not exist without them.



BY HARRY GOLDSTEIN

This replica of the original point-contact transistor is on display outside *IEEE Spectrum's* conference rooms.

Of course, this wouldn't be a Spectrum special issue if we didn't tell you how the original point-contact transistor worked, something that even the inventors seemed a little fuzzy on. According to our editorial director for content development, Glenn Zorpette, the best explanation of the point-contact transistor is in Bardeen's 1956 Nobel Prizz lecture, but even that left out important details, which Zorpette explores in classic Spectrum style in "The First Transistor and How It Worked," on page 24.

And while we're celebrating this historic accomplishment, Senior Editor Samuel K. Moore, who covers semiconductors for *Spectrum* and curated this special issue, looks at what the transistor might be like when it turns 100. For "The Transistor of 2047" [p. 38], Moore talked to the leading lights of semiconductor engineering, many of them IEEE Fellows, to get a glimpse of a druture where transistors are stacked on top of each other and are made of increasingly exotic 2D materials, even as the OG of transistor materials, germanium, is poised for a comeback.

When I was talking to Moore a few weeks ago about this issue, he mentioned that he's attending his favorite conference just as this issue comes out, the 68th edition of IEEE's International Electron Devices Meeting, in San Francisco. The mind-bending advances that emerge from that conference always get him excited about the engineering feats occurring in today's labs and on tomorrow's production lines. This year he's most excited about new devices that combine computing capability with memory to speed machine learning. Who knows, maybe the transistor of 2047 will make its debut there, too.

#### IEEE Spectrum, December 2022

### Moore's Law

Moore's law is the observation that the • number of transistors in a dense integrated circuit doubles about every two years. The observation is named after Gordon Moore, the co-founder of Fairchild Semiconductor and Intel, whose 1965 paper described a doubling every year in the number of components per integrated circuit, and projected this rate of growth would continue for at least another decade. In 1975, looking forward to the next decade, he revised the forecast to doubling every two years. The period is often quoted as 18 months because of a prediction by Intel executive David House (being a combination of the effect of more transistors and the transistors being faster).

#### WIKIPEDIA





VISUALIZING PROGRESS

# If transistors were peopl If the transistors in a microprocessor were represented by people, the following timeline gives an idea of the pace of Moore's Law.

2.300 Average music hall capacity



134,000 Large stadium capacity 32 Million

Population of Tokyo



1.3 Billion Population of China

1970	1980	1990	2000	2011
Intel 4004	Intel 286		Pentium III	Core i7 Extreme Edition

Now imagine that those 1.3 billion people could fit onstage in the original music hall. That's the scale of Moore's Law.



### Moore's Law





https://www.intel.com/content/www/us/en/newsroom/opinion/moore-law-now-and-in-the-future.html

### Transistor

• In an nMOS transistor, when a positive voltage is applied to the gate terminal relative to the source terminal, it creates an electric field that attracts electrons towards the gate. This forms a conductive channel between the source and drain terminals allowing current to flow through.



### Transistor and Gate





### Clock rate F is mainly determined by

- Switching speed of gates (transistors)
- The number of levels of gates
  - The maximum number of gates cascaded in series in any combinational logics.
  - In this example, the number of levels of gates is 3.
- Wiring delay and fanout



## Sample circuit 1 and Verilog HDL code 1

• AND gate



#### Truth table of AND gate

w_a	w_b	W_C
0	0	0
0	1	0
1	0	0
1	1	1



### Sample circuit 1 and Verilog HDL code 2

```
module top();
    wire w_c;
    reg r_a = 0;
    reg r_b = 0;
    initial #10 r_a = 1;
    initial #20 r_b = 1;
    always #1 $display("%d: %d %d %d", $time, r_a, r_b, w_c);
    initial #30 $finish();
    addgate m1(r_a, r_b, w_c);
endmodule
module addgate (
    input wire w_a,
    input wire w_b,
    output wire w_c
);
    assign #5 w_c = w_a & w_b;
endmodule
```

1:	0	0	Х	
2:	0	0	х	
3:	0	0	х	
4:	0	0	х	
5:	0	0	0	
6:	0	0	0	
7:	0	0	0	
8:	0	0	0	
9:	0	0	0	
10:	1	0	0	
11:	1	0	0	
12:	1	0	0	
13:	1	0	0	
14:	1	0	0	
15:	1	0	0	
16:	1	0	0	
17:	1	0	0	
18:	1	0	0	
19:	1	0	0	
20:	1	1	0	
21:	1	1	0	
22:	1	1	0	
23:	1	1	0	
24:	1	1	0	
25:	1	1	1	
26:	1	1	1	
27:	1	1	1	
28:	1	1	1	
29:	1	1	1	
30:	1	1	1	



## Growth in clock rate F of microprocessors

#### Intel 4004 clocked at 740KHz in 1971



#### 13th Generation Intel® Core™ i9 Processors

Products formerly Raptor Lake

Desktop	
i9-13900K	
Launched	
Q4'22	
Intel 7	
\$589.00 - \$599.00	
CPU Specifications	
Total Cores 🔞	24
# of Performance-cores	8
# of Efficient-cores	16
Total Threads 🗿	32
Max Turbo Frequency 💿	5.80 GHz
Intel® Thermal Velocity Boost Frequency ③	5.80 GHz
Intel® Turbo Boost Max Technology 3.0 Frequency ‡ 🔞	5.70 GHz
Performance-core Max Turbo Frequency 💿	5.40 GHz
Efficient-core Max Turbo Frequency 🔞	4.30 GHz
Performance-core Base Frequency	3.00 GHz
Efficient-core Base Frequency	2.20 GHz

# Syllabus (3/3)

Course	Course schedule/Required learning					
	Course schedule	Required learning				
Class 1	Design and Analysis of Computer Systems	Understand the basic of design and analysis of computer systems.				
Class 2	Instruction Set Architecture	Understand the examples of instruction set architectures				
Class 3	Memory Hierarchy Design	Understand the organization of memory hierarchy designs				
Class 4	Pipelining	Understand the idea and organization of pipelining				
Class 5	Instruction Level Parallelism: Concepts and Challenges	Understand the idea and requirements for exploiting instruction level parallelism				
Class 6	Instruction Level Parallelism: Instruction Fetch and Branch Prediction	Understand the organization of instruction fetch and branch predictions to exploit instruction level parallelism				
Class 7	Instruction Level Parallelism: Advanced Techniques for Branch Prediction	Understand the advanced techniques for branch prediction to exploit instruction level parallelism				
Class 8	Instruction Level Parallelism: Dynamic Scheduling	Understand the dynamic scheduling to exploit instruction level parallelism				
Class 9	Instruction Level Parallelism: Exploiting ILP Using Multiple Issue and Speculation	Understand the multiple issue mechanism and speculation to exploit instruction level parallelism				
Class 10	Instruction Level Parallelism: Out-of-order Execution and Multithreading	Understand the out-of-order execution and multithreading to exploit instruction level parallelism				



### From multi-core era to many-core era



Figure 1: Current and expected eras of Intel® processor architectures

#### Platform 2015: Intel® Processor and Platform Evolution for the Next Decade, 2005

### Pollack's Rule

 Pollack's Rule states that microprocessor performance increase due to microarchitecture advances is roughly proportional to the square root of the increase in complexity. Complexity in this context means processor logic, i.e. its area.



WIKIPFDIA

### From multi-core era to many-core era



EV6	EV6	EV6
EV6	EV6	EV6
EV6	EV6	EV6

#### Figure 1. Relative sizes of the cores used in the study

Single-ISA Heterogeneous Multi-Core Architectures: The Potential for Processor Power Reduction, MICRO-36

## Intel Sandy Bridge, January 2011

• 4 core

Processor	Core	Core	Core	Core	System Agent & Memory
Graphics a					Controller including
		Shared L3	Cache**		and Misc. I/O
- HELL - AND - CONT		Memory	Controller I/O		



## Intel Skylake-X, Core i9-7980XE, 2017

• 18 core







# Syllabus (3/3)

#### Course schedule/Required learning Course schedule Required learning Design and Analysis of Computer Systems Understand the basic of design and analysis of Class 1 computer systems. Class 2 Instruction Set Architecture Understand the examples of instruction set architectures Class 3 Memory Hierarchy Design Understand the organization of memory hierarchy designs Class 4 Pipelining Understand the idea and organization of pipelining Class 5 Instruction Level Parallelism: Concepts and Challenges Understand the idea and requirements for exploiting instruction level parallelism Understand the organization of instruction fetch and Instruction Level Parallelism: Instruction Fetch and Branch Prediction Class 6 branch predictions to exploit instruction level parallelism Class 7 Instruction Level Parallelism: Advanced Techniques for Branch Prediction Understand the advanced techniques for branch prediction to exploit instruction level parallelism Instruction Level Parallelism: Dynamic Scheduling Understand the dynamic scheduling to exploit Class 8 instruction level parallelism Instruction Level Parallelism: Exploiting ILP Using Multiple Issue and Class 9 Understand the multiple issue mechanism and Speculation speculation to exploit instruction level parallelism Class 10 Instruction Level Parallelism: Out-of-order Execution and Multithreading Understand the out-of-order execution and multithreading to exploit instruction level parallelism Class 11 Multi-Processor: Distributed Memory and Shared Memory Architecture Understand the distributed memory and shared memory architecture for multi-processors Class 12 Thread Level Parallelism: Coherence and Synchronization Understand the coherence and synchronization for thread level parallelism Understand the memory consistency model for thread Class 13 Thread Level Parallelism: Memory Consistency Model level parallelism Thread Level Parallelism: Interconnection Network and Man-core Class 14 Understand the interconnection network and many-Processors core processors for thread level parallelism



## Adaptive Computing Research Initiative (ACRi)

- The aim
  - Aiming to develop the high-performance Adaptive Computing Systems that utilize FPGAs
  - Working out to distribute the FPGA-related technologies, including our developed systems, as an outreach activity for research results
- Main research theme
  - 1. Development for FPGA accelerator to speed up processing of AI etc.
  - 2. Development for FPGA accelerators and FPGA systems for **IoT**.
- Activity
  - Establishment Date: April 1<sup>st</sup> 2020
  - Activity period : First period of 3 years



The Adaptive Computing Research Initiative is an organization to seek out and research ways to utilize FPGAs.

# Please apply for your user account on this site today

https://gw.acri.c.titech.ac.jp/wp/manual/apply-for-account

ACRi	
ACRi ルームのアカウント申請方法	サイト内を検索 Q
© 2020.08.11 © 2020.06.30	ACRi ルームの情報
□ 目次 [閉じる]	
	ようこそ
<ol> <li>アカリントの申請</li> <li>登録フォームへの入力</li> </ol>	予約ページトップ
<ul> <li>アカウント申請の受理</li> </ul>	ニュースとメンテナンス情報
2. ログインおよびパスフレーズの変更	フォーラム
<ul> <li>予約システム</li> </ul>	ギャラリーと技術情報
・ ACRi のサーバ	
アカウントの申請	ログイン/ログアウト
	ログイン
登録フォームへの入力	
申請するには、 <u>「アカウントの申請」のリンク</u> をクリックするか、ログインが必要なページ 例えば各サーバの予約状況のページ) にアクセスしてください。後者の場合には、ログイン フォームの右下にある「新規ユーザー登録」をクリックしてください。	ACRiルームの利用説明



### Discussion: software and hardware

```
#include <stdio.h>
main()
{
    printf("hello, world¥n");
}
```



Hardware to light up some LEDs

