Clock, Power Consumption and the Future Landscape of Computation

Prof. Usagi

Recap: Moore's Law

 The number of transistors we can build in a fixed area of silicon doubles every 12 ~ 24 months.





2015

Recap: Pipelining a 4-bit serial adder

add	а,	b
add	С,	d
add	e,	f
add	g,	h
add	i,	j
add	k,	1
add	m,	n
add	Ο,	р
add	q,	r
add	S,	t
add	u,	V

								C	alaa				
1st	2nd	3rd	4th					Cyc	les		1		
	1st	2nd	3rd	4th				A	47	·			
		1st	2nd	3rd	4th			110	лU				
			1st	2nd	3rd	4th							
				1st	2nd	3rd	4th						
					1st	2nd	3rd	4th					
						1st	2nd	3rd	4th				
					•		1st	2nd	3rd	4th			
				Afte	er thi	s poi	nt,	1st	2nd	3rd	4th		
					are c	omp		g an	1st	2nd	3rd	4th	
					lope	ratio	nea	JU		1st	2nd	3rd	4th
			-										



t

Recap: The growth of clock rate is slowing down



el Skylake Core i7 00 MHz in 2017						
2%/year						
0 2012 2014 2016 2018						

Outline

- What are the basic limits of clock frequency?
- New limit on clock frequency: Power consumption
- Opportunities and the future

Timing constraints



Output Timing Constraints



- Min delay of FF, also called contamination delay or min CLK to Q delay: t_{cca}
 - Time after clock edge that Q might be unstable (i.e., starts changing)
- Max delay of FF, also called propagation delay or maximum CLK to Q delay: t_{pcq}
 - Time after clock edge that the output Q is guaranteed to be stable (i.e. stops) changing)



 t_{pcq}

Setup and hold times for a flip-flop



- Setup time: *t_{setup}*
 - Time before the clock edge that data must be stable (i.e. not change)
 - For the FF to capture the input
- Hold time: *t*_{hold}
 - Time after the clock edge that data must be stable
 - For the FF to output/store/propagate the data
- Aperture time: t_a
 - Time around clock edge that data must be stable ($t_a = t_{setup} + t_{hold}$)





Summary on timing constraints

- Combinational:
 - Maximum delay = **P**ropagation **d**elay (t_{pd})
 - Minimum delay = **C**ontamination **d**elay (t_{cd})
- Flip Flops:
 - Input

- **Once the logic/FFs are built,** these timing characteristics are fixed properties
- **Setup** time (*t_{setup}*)
- Hold time (*t*_{hold})
- Output
 - **P**ropagation **c**lock-to-**Q** time (t_{pcq})
 - Contamination clock-to-Q time (t_{ccq})



D1

R1



Timing in a circuit



Causes of Timing Issues in Sequential Circuits

- Input to a FF comes from the output of another FF through a combinational circuit
- The FF and combinational circuit have a min & max delay
- Which of the following violations occurs if max delay of R1 is zero & max delay of the combinational circuit is equal to the clock period?
 - A. Hold time violation for R2
 - B. Setup violation for R2
 - C. Hold time violation for R1
 - D. Setup violation for R1
 - E. None of the above

R2

Combinational Logic

R1

C[K

Causes of Timing Issues in Sequential Circuits

- Input to a FF comes from the output of another FF through a combinational circuit
- The FF and combinational circuit have a min & max delay
- Which of the following violations occurs if max delay of R1 is zero & max delay of the combinational circuit is equal to the clock period?
 - A. Hold time violation for R2
 - B. Setup violation for R2
 - C. Hold time violation for R1
 - D. Setup violation for R1
 - E. None of the above

 $T_{c} \geq t_{setup} + \max_{delay(FF)} + \max_{delay(combinational)} \boxed{Q1}$ $T_{c} \geq t_{setup} + t_{pcq} + t_{pd} = t_{c}$ $T_{c} \geq t_{setup} + t_{pcq} + t_{pd} = t_{c}$

R1

CLK

Combinational Logic	→ R2
T_	
	/
	I
	I
	i

Poll close in 1:30

Causes of Timing Issues in Sequential Circuits (2)

- Input to a FF comes from the output of another FF through a combinational circuit
- The FF and combinational circuit have a min & max delay Which of the following violations occurs if min delay of R1 is zero & max delay of the combinational circuit was just a wire?
 - A. Hold time violation for R2
 - B. Setup violation for R2
 - C. Hold time violation for R1
 - D. Setup violation for R1
 - E. None of the above

Combinational Logic

R1

CĹK

R2

Causes of Timing Issues in Sequential Circuits (2)

- Input to a FF comes from the output of another FF through a combinational circuit
- The FF and combinational circuit have a min & max delay Which of the following violations occurs if min delay of R1 is zero & max delay of the combinational circuit was just a wire?

R1

CLK

A. Hold time violation for R2

- B. Setup violation for R2
- C. Hold time violation for R1
- D. Setup violation for R1

E. None of the above $T_{hold} \leq \min_delay(FF) + \min_delay(combinational)^{-1}$ Q1 $T_{hold} \leq t_{cca} + t_{cl}$ D2 14

-	.
Combinational Logic	R2
–	́Ц
۱ _с	
	/
	· ·



Setup time constraints $T_c \ge t_{pcq} + t_{pd} + t_{setup} + t_{skew}$ Hold time constraints $t_{ccq} + t_{cd} > t_{hold}$ $30ps + 25ps > t_{hold}$

 $t_{hold} = 70 \text{ ps!}$ **No!!!**

Flip flops						
t _{ccq}	30 ps					
<i>t</i> _{pcq}	50 ps					
t _{setup}	60 ps					
thold	70 ps					



Setup time constraints

Hold time constraints $t_{ccq} + t_{cd} > t_{hold}$

Flip flops

<i>t_{ccq}</i>	30 ps
<i>t</i> _{pcq}	50 ps
t setup	60 ps
thold	70 ps



Example: timing constraints

- What's the maximum frequency?
 - A. 1/110ns
 - B. 1/220ns
 - C. 1/200ns
 - D. 1/180ns
 - E. None of the above

$$T_{c} \geq t_{pcq} + t_{pd} + t_{setup}$$
$$t_{ccq} + t_{cd} > t_{hold}$$





Example: timing constraints

- What's the maximum frequency?
 - A. 1/110ns

B. 1/220ns

- C. 1/200ns
- D. 1/180ns
- E. None of the above

 $T_c \ge t_{pcq} + t_{pd} + t_{setup}$ $T_c \ge 70ns + 130ns + 20ns + 0$





FF Timing Parameters

- Once a flip flop has been built, its timing characteristics stay fixed: *t_{setup}*, *t_{hold}*, *t_{ccq}*, *t_{pcq}*
- What about the clock? Does the clock edge arrive at the same time to all the D-FFs on the chip?









解除靜音 取用 **Minimum** Minimum number of SC function?

- A. 1
- B. 2
- C. 3
- D. 4 E. 5

Clock Skew

- The clock doesn't arrive at all registers at the same time
- Skew: difference between the two clock edges
- Perform the worst case analysis



The wire has its own delay!!!

Setup Time Constraint with Skew

- In the worst case, CLK2 is earlier than CLK1
- t_{pcq} is max delay through FF, t_{pd} is max delay through logic





$T_c \ge t_{pcq} + t_{pd} + t_{setup} + t_{skew}$ $t_{pd} \le T_c - (t_{setup} + t_{pcq} + t_{skew})$

The larger the design, the longer the t_{skew}

Power consumption



Power & Energy

- Regarding power and energy, how many of the following statements are correct?
 - ① Lowering the power consumption helps extending the battery life
 - Lowering the power consumption helps reducing the heat generation (2)
 - ③ Lowering the energy consumption helps reducing the electricity bill
 - ④ A CPU with 10% utilization can still consume 33% of the peak power
 - A. 0
 - B. 1
 - C. 2

D. 3

E. 4

Power v.s. Energy

- Power is the direct contributor of "heat"
 - Packaging of the chip
 - Heat dissipation cost
 - Power = $P_{Dynamic} + P_{static}$
- Energy = P * ET
 - The electricity bill and battery life is related to energy!
 - Lower power does not necessary means better battery life if the processor slow down the application too much

ergy! Pattery life if the

Power & Energy

 Regarding power and energy, how many of the following statements are correct?

① Lowering the power consumption helps extending the battery life Lowering the power consumption helps reducing the heat generation Lowering the energy consumption helps reducing the electricity bill ④ A CPU with 10% utilization can still consume 33% of the peak power

- A. 0
- **B**. 1
- C. 2

D. 3

F. 4

Dynamic/Active Power

- The power consumption due to the switching of transistor states
- Dynamic power per transistor $P_{dynamic} \sim \alpha \times C \times V^2 \times f \times N$
 - α : average switches per cycle
 - C: capacitance
 - V: voltage
 - f: frequency, usually linear with V
 - N: the number of transistors



Static/Leakage Power

- The power consumption due to leakage transistors do not turn all the way off during no operation
- Becomes the dominant factor in the most advanced process technologies. 1000

$$P_{leakage} \sim N \times V \times e^{-V_t}$$

- N: number of transistors
- V: voltage
- V_t : threshold voltage where transistor conducts (begins to switch)



Figure 1: Leakage power becomes a growing problem as demands for more performance and functionality drive chipmakers to nanometer-scale process nodes (Source: IBS).



Power & Energy

 Regarding power and energy, how many of the following statements are correct?

① Lowering the power consumption helps extending the battery life Lowering the power consumption helps reducing the heat generation Lowering the energy consumption helps reducing the electricity bill A CPU with 10% utilization can still consume 33% of the peak power A. 0

- **B**. 1
- C. 2

D. 3

Dynamic/Active Power

- The power consumption due to the switching of transistor states
- Dynamic power per transistor $P_{dynamic} \sim \alpha \times C \times V^2 \times f \times N$
 - α : average switches per cycle
 - C: capacitance
 - *V*: voltage
 - f: frequency, usually linear with V
 - N: the number of transistors



Dennardian Broken

Given a scaling factor S

Parameter	Relation	Classical Scaling	Leakage Limited
Power Budget		1	1
Chip Size		1	1
Vdd (Supply Voltage)		1/S	1
Vt (Threshold Voltage)	1/S	1/S	1
tex (oxide thickness)		1/S	1/S
W, L (transistor dimensions)		1/S	1/S
Cgate (gate capacitance)	WL/tox	1/S	1/S
Isat (saturation current)	WVdd/tox	1/S	1
F (device frequency)	Isat/(CgateVdd)	S	S
D (Device/Area)	1/(WL)	S ²	S ²
p (device power)	IsatVdd	1/S ²	1
P (chip power)	Dp	1	S ²
U (utilization)	1/P	1	1/S ²

Poll close in 1:30

What happens if power doesn't scale with process technologies?

- If we are able to cram more transistors within the same chip area (Moore's law continues), but the power consumption per transistor remains the same. Right now, if put more transistors in the same area because the technology allows us to. How many of the following statements are true?
 - ① The power consumption per chip will increase
 - ² The power density of the chip will increase
 - Given the same power budget, we may not able to power on all chip area if we maintain the (3) same clock rate
 - Given the same power budget, we may have to lower the clock rate of circuits to power on all (4) chip area
 - A. 0
 - B. 1
 - C. 2
 - D. 3
 - E. 4

Power consumption

Chip									
1	1	1	1	1	1	1			
1	1	1	1	1	1	1			
1	1	1	1	1	1	1			
1	1	1	1	1	1	1			
1	1	1	1	1	1	1			
1	1	1	1	1	1	1			
1	1	1	1	1	1	1			

Dennardian Scaling

Chip

0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5

=50W



Dennardian Broken

Chip										
1	1	1	1	1	1	1	1	1		
1	1	1	1	1	1	1	1	1		
1	1	1	1	1	1	1	1	1		
1	1	1	1	1	1	1	1	1		
1	1	1	1	1	1	1	1	1		
1	1	1	1	1	1	1	1	1		
1	1	1	1	1	1	1	1	1		
1	1	1	1	1	1	1	1	1		
1	1	1	1	1	1	1	1	1		
1	1	1	1	1	1	1	1	1		

=100W!

Power density

Chip

Dennardian Scaling

Chip

0.5



50W Chip Area

Jennardian Broken											
Chip											
1	1	1	1	1	1	1	1	1			
1	1	1	1	1	1	1	1	1			
1	1	1	1	1	1	1	1	1			
1	1	1	1	1	1	1	1	1			
1	1	1	1	1	1	1	1	1			
1	1	1	1	1	1	1	1	1			
1	1	1	1	1	1	1	1	1			
1	1	1	1	1	1	1	1	1			
1	1	1	1	1	1	1	1	1			
1	1	1	1	1	1	1	1	1			

100W

Chip Area

Power density



Figure 1. In CPU architecture today, heat is becoming an unmanageable problem. (Courtesy of Pat Gelsinger, Intel Developer Forum, Spring 2004)

What happens if power doesn't scale with process technologies?

- If we are able to cram more transistors within the same chip area (Moore's law continues), but the power consumption per transistor remains the same. Right now, if put more transistors in the same area because the technology allows us to. How many of the following statements are true?
 - The power consumption per chip will increase
 - The power density of the chip will increase
 - ③ Given the same power budget, we may not able to power on all chip area if we maintain the same clock rate
 - Given the same power budget, we may have to lower the clock rate of circuits to power on all (4) chip area
 - A. 0
 - B. 1
 - C. 2
 - D. 3
 - E. 4

Power consumption to light on all transistors If we can only cool down 50W in the same area — Dennardian Scaling Dennardian Broken

Chip												
1	1	1	1	1	1	1						
1	1	1	1	1	1	1						
1	1	1	1	1	1	1						
1	1	1	1	1	1	1						
1	1	1	1	1	1	1						
1	1	1	1	1	1	1						
1	1	1	1	1	1	1						

Chip

0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5
0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5	0.5

=50W

=49W



=100W!

Dark silicon

- Your power consumption goes up as the number of transistors goes up
- Even Moore's Law allows us to put more transistors within the same area —we cannot use them all simultaneously!
- We have no choice to not activate all transistors at the same time!

What happens if power doesn't scale with process technologies?

- If we are able to cram more transistors within the same chip area (Moore's law continues), but the power consumption per transistor remains the same. Right now, if put more transistors in the same area because the technology allows us to. How many of the following statements are true?
 - The power consumption per chip will increase
 - The power density of the chip will increase
 - Given the same power budget, we may not able to power on all chip area if we maintain the same clock rate
 - Given the same power budget, we may have to lower the clock rate of circuits to power on all (4) chip area
 - A. 0
 - B. 1
 - C. 2
 - D. 3
 - E. 4

Dynamic/Active Power

- The power consumption due to the switching of transistor states
- Dynamic power per transistor $P_{dynamic} \sim \alpha \times C \times V^2 \times f \times N$
 - α : average switches per cycle
 - C: capacitance
 - *V*: voltage
 - f: frequency, usually linear with V
 - N: the number of transistors



What happens if power doesn't scale with process technologies?

- If we are able to cram more transistors within the same chip area (Moore's law continues), but the power consumption per transistor remains the same. Right now, if put more transistors in the same area because the technology allows us to. How many of the following statements are true?
 - **W** The power consumption per chip will increase
 - The power density of the chip will increase
 - Given the same power budget, we may not able to power on all chip area if we maintain the same clock rate
 - Given the same power budget, we may have to lower the clock rate of circuits to power on all chip area
 - A. 0
 - B. 1
 - C. 2
 - D. 3

Solutions/trends in dark silicon era

Trends in the Dark Silicon Era

- Aggressive dynamic voltage/frequency scaling
- Throughout oriented slower, but more
- Just let it dark activate part of circuits, but not all
- From general-purpose to domain-specific ASIC



Aggressive dynamic frequency scaling

Dynamic/Active Power

- The power consumption due to the switching of transistor states
- Dynamic power per transistor $P_{dynamic} \sim \alpha \times C \times V^2 \times f \times N$
 - α : average switches per cycle
 - C: capacitance
 - *V*: voltage
 - f: frequency, usually linear with V
 - N: the number of transistors



Frequency varies per core

Products	Solutions Support		(intel)
		X Intel® Xeon® Processor E7-8890 v4	Intel® Xeon® Processor E7-8893 v4
	Status	Launched	Launched
	Launch Date 🚯	Q2"16	Q2'16
	Lithography 🕄	14 nm	14 nm
	Performance		
	# of Cores 🜖	24	4
	# of Threads 🟮	48	8
	Processor Base Frequency 🧕	2.20 GHz	3.20 GHz
	Max Turbo Frequency 🚺	3.40 GHz	3.50 GHz
	Cache 🚯	60 MB	60 MB
	Bus Speed 🚯	9.6 GT/s	9.6 GT/s
	# of QPI Links 🚯	3	3
	TDP 🟮	165 W	140 W

. . . .



×	Intel® Xeon® Processor E7–8880 v4	×	
	Launched		
	Q2'16		
	14 nm		
	22		
	44		
	2.20 GHz		
	3.30 GHz		
	55 MB		

9.6 GT/s

З

150 W

Demo

- You may use cat /proc/cpuinfo to see all the details of your processors
- You may add "| grep MHz" to see the frequencies of your cores
- Only very few of them are on the boosted frequency

Slower, but more



More cores per chip, slower per core

Products	Solutions Support		(intel)
		X Intel® Xeon® Processor E7-8890 v4	Intel® Xeon® Processor E7-8893 v4
	Status	Launched	Launched
	Launch Date 🧕	Q2'16	Q2'16
	Lithography 🟮	14 nm	14 nm
	Performance		
	# of Cores 🜖	24	4
	# of Threads 🟮	48	8
	Processor Base Frequency 🧿	2.20 GHz	3.20 GHz
	Max Turbo Frequency 🕕	3.40 GHz	3.50 GHz
	Cache 🚯	60 MB	60 MB
	Bus Speed 🚯	9.6 GT/s	9.6 GT/s
	# of QPI Links 🟮	3	3
	TDP 🟮	165 W	140 W

.

×	Intel® Xeon® Processor E7-8880 v4	×	
	Launched		
	Q2'16		
	14 nm		
	22		
	44		
	2.20 GHz		
	3.30 GHz		

55 MB

9.6 GT/s

3

150 W



Each of these performs the same operation, but each of these is also a "thread"

SMX	5) 5)							laa	tru oti	on Co	aha							
	War	n Sch	eduler	-	-	Wa	rn Scher	Ins	urucu	on ca	Che Wai	n Set	odulor	_		Wa	rn Schor	lulor
Dispatch Dispatch			Dispatch Dispatch						Dispatch Dispatch					Dispatch Dispa				
	Ŧ		Ŧ			Ŧ		÷			Ŧ		Ŧ			Ŧ		Ŧ
				Regi	ster	File (65,536	x 32-	bit G	K110) (1:	31,07	2 x32-k	oit Gl	<210)			
	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+	+
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	OF Jni	0	ld	ore	P_Jnit	.D. G	SFU		Cole	ore		re	C		25	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST
	1	64 KI	B Shar	ed Ma	emor	v / L 4	1 Cache		110) I	(128	KBS	Share	ed Mem	ory	110	ache	GK21))
	V		- Active			, <u> </u>	48 K	B Re	ad-O	nly D	ata C	ache)	ery i			GINE	0
	Tex		Тех			Tex		Тех	(Tex		Tex	(Tex		Tex
	Tex		Тех	(Tex		Tex	¢		Tex		Tex	(Tex		Tex
									ť	DI								



ARM's big.LITTLE architecture big.LITTLE system





Just let it dark

NVIDIA's Turing Architecture







Programming in Turing Architecture

Use tensor cores cublasErrCheck(cublasSetMathMode(cublasHandle, CUBLAS_TENSOR_OP_MATH)); Make them 16-bit convertFp32ToFp16 <<< (MATRIX_M * MATRIX_K + 255) / 256, 256 >>> (a_fp16, a_fp32, $MATRIX_M * MATRIX_K);$ convertFp32ToFp16 <<< (MATRIX_K * MATRIX_N + 255) / 256, 256 >>> (b_fp16, b_fp32, MATRIX_K * MATRIX_N);

cublasErrCheck(cublasGemmEx(cublasHandle, CUBLAS_OP_N, CUBLAS_OP_N, MATRIX_M, MATRIX_N, MATRIX_K, &alpha, a_fp16, CUDA_R_16F, MATRIX_M, b_fp16, CUDA_R_16F, MATRIX_K, &beta, c_cublas, CUDA_R_32F, MATRIX_M, CUDA_R_32F, CUBLAS_GEMM_DFALT_TENSOR_OP));

call Gemm

NVIDIA's Turing Architecture





The rise of ASICs





Real-time AI: Microsoft announces preview of Project Brainwave



Blog Latest Stories Product News Topics

AI & MACHINE LEARNING

An in-depth look at Google's first Tensor Processing Unit (TPU)

Kaz Sato

Staff Developer Advocate, Google Cloud There's a common thread that connects Google services such as Google Search, Street View, Google Photos and Google Translate: they all use Google's Tensor Processing Unit, or TPU, to accelerate their neural network computations behind the scenes.

Cliff Young Software Engineer, Google Brain

David Patterson Distinguished Engineer, Google Brain

May 12, 2017

Try GCP

Get \$300 free credit to spend over 12 months.



Google's first Tensor Processing Unit (TPU) on a printed circuit board (left); TPUs deployed in a Google datacenter (right)



There is no pure "software" or "hardware" design in the dark silicon era. Everything needs to be hardware/software co-designed.

–Prof. Usagi

Announcement

- iEval Capture your screenshot, submit through iLearn and you will receive a full credit assignment
- Assignment 6 due 6/4
- Lab 6 due this Friday
- Please fill out ABET survey through iLearn
- Final exam will be held during the campus scheduled period to avoid conflicts
 - Final review -6/4 during the lecture, will also release the sample final
 - 6/11 11:30am 2:59:59pm
 - About the same format as midterm, but longer
 - Will have a final review on 6/6 to help you prepare

Electrical Computer Science Engineering





