Dark Silicon & Modern Computer Architecture

Hung-Wei Tseng

Recap: Take-aways of parallel programming

- Processor behaviors are non-deterministic
 - You cannot predict which processor is going faster
 - You cannot predict when OS is going to schedule your thread
- Cache coherency only guarantees that everyone would eventually have a coherent view of data, but not when
- Cache consistency is hard to support





Dark Silicon and its impact on computer architecture

Power and Energy



Power & Energy

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



https://www.pollev.com/hungweitseng close in 1:30

Power&Energy

А

С

D

Ε

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

Dynamic Power

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



Double Clock Rate or Double the # of Processors?

 Assume 60% of the application can be fully parallelized with 2core or speedup linearly with clock rate. Should we double the **clock rate or duplicate a core?**

$$P_{dynamic} \sim \alpha \times C \times V^2 \times f \times N$$

 $Speedup_{parallel}(f_{parallelizable}, n) = \frac{1}{(1 - f_{parallelizable}) + \frac{f_{parallelizable}}{n}}$ Speedup_{parallel}(60 %,2) = $\frac{1}{(1-60\%) + \frac{60\%}{2}} = 1.43$

 $Power_{2-core} = 2 \times P_{baseline}$ $Speedup_{2 \times clock} = 2$ $Energy_{2-core} = 2 \times P_{baseline} \times ET_{baseline} \times \frac{1}{1.43} = 1.39 \times Energy_{baseline}$ $Power_{2 \times clock} = 2^3 \times P_{baseline} = 8 \times P_{baseline}$

 $Energy_{2 \times clock} = 2^3 \times P_{baseline} \times ET_{baseline} \times \frac{1}{2} = 4 \times P_{baseline} \times ET_{baseline}$

Dynamic voltage/frequency scaling

- Dynamically lower power for performance
 - Change the voltage and frequency at runtime
 - Under control of operating system that's why updating iOS may slow down an old iPhone
- Recall: $P_{dynamic} \sim \alpha \times C \times V^2 \times f \times N$
 - Because frequency ~ to V...
 - $P_{dynamic} \sim to V^3$
- Reduce both V and f linearly
 - Cubic decrease in dynamic power
 - Linear decrease in performance (actually sub-linear)
 - Thus, only about quadratic in energy
 - Linear decrease in static power
 - Thus, only modest static energy improvement
 - Newer chips can do this on a per-core basis
 - cat /proc/cpuinfo in linux



Demo — changing the max frequency and performance

- Change the maximum frequency of the intel processor you learned how to do this when we discuss programmer's impact on performance
- LIKWID a profiling tool providing power/energy information
 - likwid-perfctr -g ENERGY [command_line]
 - Let's try blockmm and popcount and see what's happening!

Power & Energy

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



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
 - The power density of the chip will increase 2
 - 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

Α.	0	
B.	1	
C.	2	
D.	3	
E.	4	

https://www.pollev.com/hungweitseng close in 1:30

Dark Silicon

А



D

С

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

Dark Silicon and the End of Multicore Scaling

H. Esmaeilzadeh, E. Blem, R. St. Amant, K. Sankaralingam and D. Burger University of Washington, University of Wisconsin—Madison, University of Texas at Austin, Microsoft Research

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).



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)	lsat/(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

Power consumption to light on all transistors

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

=50W

=49W

Dennardian Broken



=100W!

Clock rate improvement is limited nowadays





Skylake Core i7 0 MHz in 2017	,
%/year	
201 201 2010 1	0,2

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

Modern processor's frequency



Intel[®] Core[™] i9-9900K Processor

16M Cache, up to 5.00 GHz

Essentials	CPU Specifications	
Product Collection	# of Cores 🕐	8
Code Name	# of Threads 🕜	16
Vertical Segment Processor Number (?)	Processor Base Frequency 🕐	3.60 GHz
Status	Max Turbo Frequency 🕐	5.00 GHz
Launch Date 🕐	Cache 🕐	16 MB Intel® Sma
Lithography 🕐	Bus Speed 🕐	8 GT/s
Included Items	Intel® Turbo Boost Technology 2.0 Frequency‡ 👔	5.00 GHz
Recommended Customer Price (2)	TDP 🕐	95 W





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



Recap: Demo — changing the max frequency and performance

- Change the maximum frequency of the intel processor you learned how to do this when we discuss programmer's impact on performance
- LIKWID a profiling tool providing power/energy information
 - likwid-perfctr -g ENERGY [command_line]
 - Let's try blockmm and popcount and see what's happening!

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.
 - Pleakage ~ How about static pow
 - 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).



Slower, but more



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

Rakesh Kumar, Keith I. Farkas*, Norman P. Jouppi*, Parthasarathy Ranganathan*, Dean M. Tullsen

UCSD and HP Labs*

In the 36th Annual IEEE/ACM International Symposium on Microarchitecture, 2003. MICRO-36., 2003 MICRO Test-of-time award, 2021

Areas of different processor generations

- You fit about 5 EV5 cores within the same area of an EV6
- If you build a quad-core EV6, you can use the same area to
 - build 20-core EV5
 - 3EV6+5EV5

Processor	EV5	EV6	EV6+
Issue-width	4	6 (000)	6 (000)
I-Cache	8KB, DM	64KB, 2-way	64KB, 2-way
D-Cache	8KB, DM	64KB, 2-way	64KB, 2-way
Branch Pred.	2K-gshare	hybrid 2-level	hybrid 2-level
Number of MSHRs	4	8	16
Number of threads	1	1	4
Area (in mm^2)	5.06	24.5	29.9



https://www.pollev.com/hungweitseng close in 1:30 **Single ISA heterogeneous CMP**

- Regarding "Single-ISA Heterogeneous Multi-Core Architectures", how many of the following statements is/are correct?
 - ① You need to recompile and optimize the binary for each core architecture to exploit the threadlevel parallelism in this architecture
 - 2 For a program with limited thread-level parallelism, single ISA heterogeneous CMP would deliver better or at least the same level of performance than homogeneous CMP
 - ③ For a program with rich thread-level parallelism, single ISA heterogeneous CMP would deliver better or at least the same level of performance than homogeneous CMP built with oldergeneration cores
 - ④ Spending more instructions on older-generation cores would always lead to better energy-delay **Big.Little**
 - A. 0
 - B. 1
 - C. 2
 - D. 3
 - E. 4

Ε

Energy-delay

Energy * delay = Power * ET * ET = Power * ET²



Benchmark	Total	%	of instru	ctions per	core	Energy	ED	ED^2	Perf.
	switches	EV4	EV5	EV6	EV8-	Savings(%)	Savings(%)	Savings(%)	Loss (%)
ammp	0	0	0	0	100	0	0	0	0
applu	27	2.2	0.1	54.5	43.2	42.7	38.6	33.6	7.1
apsi	2	0	0	62.2	37.8	27.6	25.3	22.9	3.1
art	0	0	0	100	0	74.4	73.5	72.6	3.3
equake	20	0	0	97.9	2.1	72.4	71.3	70.1	3.9
fma3d	0	0	0	0	100	0	0	0	0
wupwise	16	0	0	99	1	72.6	69.9	66.2	10.0
bzip	13	0	0.1	84.0	15.9	40.1	38.7	37.2	2.3
crafty	0	0	0	0	100	0	0	0	0
eon	0	0	0	100	0	77.3	76.3	75.3	4.2
gzip	82	0	0	95.9	4.1	74.0	73.0	71.8	3.9
mcf	0	0	0	0	100	0	0	0	0
twolf	0	0	0	0	100	0	0	0	0
vortex	364	0	0	73.8	26.2	56.2	51.9	46.2	9.8
Average	1(median)	0.2%	0%	54.8%	45.0%	38.5%	37.0%	35.4%	3.4%

4EV6 v.s. 20 EV5 v.s. 3EV6+5EV5







4EV6 v.s. 20 EV5 v.s. 3EV6+5EV5





Single ISA heterogeneous CMP

- Regarding "Single-ISA Heterogeneous Multi-Core Architectures", how many of the following statements is/are correct?
 - ① You need to recompile and optimize the binary for each core architecture to exploit the thread-level parallelism in this architecture
 - For a program with limited thread-level parallelism, single ISA heterogeneous CMP would deliver better or at least the same level (less than 10% difference) of performance than homogeneous CMP
 - ③ For a program with rich thread-level parallelism, single ISA heterogeneous CMP would deliver better or at least the same level of performance than homogeneous CMP built with oldergeneration cores
 - ④ Spending more instructions on older-generation cores would always lead to better energy-delay





ARM's big.LITTLE architecture big.LITTLE system





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 1	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

Xeon Phi

Essentials

Product Collection	Intel® Xeon Phi [™] 72x5 Processor Family
Code Name	Products formerly Knights Mill
Vertical Segment	Server
Processor Number	7295
Off Roadmap	No
Off Roadmap Status	No Launched
Off Roadmap Status Launch Date 🕐	No Launched Q4'17

Performance

# of Cores 👔	72
of Threads (72
Processor Base Frequency 🕜	1.50 GHz
Max Turbo Frequency 🕐	1.60 GHz
Cache 🕐	36 MB L2 Cache
TDP 👔	320 W



The rise of GPU



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

SMX	s: E							Inc	tructi	on Ca	cho								
Warp Scheduler					-	Warp Scheduler Warp Scheduler										Warn Scheduler			
Dispatch Dispatch			Dispatch Dispatch						Dispatch Dispatch				Dispatch Dispa						
	÷		÷		h.	+		÷			Ŧ		Ŧ		181 <mark>1</mark> 11270 - 2008	+		Ŧ	
				Regi	ster	File (65,536	x 32-	bit G	K110) (1:	31,07	'2 x32-k	oit Gl	< 210))			
							+										+	-	
Core	Core	Core	UP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	UP 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	DF Uni	0	51	ore	P Jnit	.D, C	SFU		Core	Core	S.	Are	CC		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	
		64 KI	B Shar	ad Ma	mor	v/14	1 Cache		110\ I	(128	KRS	Share	od Mom	ory	110	ache	GK21	n \	
1	(o onan			y ,	48 K	B Re	ad-O	nlv D	ata C	ache		or y /	EIU	achte		7	
	Tex		Tex			Tex		Тех	¢	Tex			Тех		Tex			Tex	
	Tex		Tex	:		Tex		Tex	٢		Tex		Te>	(Tex		Tex	
									Z	19									



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



Say, we want to implement a[i] += a[i+1]*20

This is what we need in RISC-V in each iteration





This is what you need for these instructions









We don't need these many registers, complex control, decode

We don't need instruction fetch given it's a fixed function













Rearranging the datapath

ld X1, 0(X0)
ld X2, 8(X0)
add X3, X31, #20
mul X2, X2, X3
add X1, X1, X2
sd X1, 0(X0)



The pipeline for a[i] += a[i+1]*20

a[3] += a[4]*20 a[2] += a[3]*20 Each stage can still be as fast as the pipelined processor Data But each stage is memory Read Data Register now working on Address what the original 6 Read instructions would Register Address Data Adder 8 Data do memory 20.



Announcement

- Project due tonight try your best
- Assignment #4 due this Wednesday
- iEVAL, until 12/3
 - Please fill the survey to let us know your opinion!
 - Don't forget to take a screenshot of your submission and submit through iLearn it counts as a **full credit assignment**
 - We will drop your lowest 2 assignment grades
- Final Exam
 - Starting from 12/6 to 12/10 12:00pm, any consecutive 180 minutes you pick
 - Similar to the midterm, but more time and about 1.5x longer
 - Two of the problem sets will be comprehensive exam questions
 - Will release a sample final at the end of the last lecture