## **Midterm Review**

Hung-Wei Tseng

# von Neumann model & Performance equation

### von Neumman Architecture



### By loading different programs into memory, your computer can perform different functions







### 00005d24 00c2f000 nstructio ta 0000bd24 0000008 0 0 2ca422a0 00c2f800 80000008 130020e4 00003d24 00c30000 2ca4e2b3 80000008

Storage

## **Execution Time**

- The simplest kind of performance
- Shorter execution time means better performance
- Usually measured in seconds



### instruction memory

```
ah gp,15(t12)

a gp,-25520(gp)

ah t1,0(gp)

ah t4,0(gp)

l t0,-23508(t1)

q t0,120007a94

ah t0,0(gp)

l zero,-23508(t1)

r v0

l zero,-23512(t4)

q t0,120007a98

v t0,t1

r t2

120007a80
```

### **CPU Performance Equation**



# Execution Time = $\frac{Instructions}{Program} \times \frac{Cycles}{Instruction} \times \frac{Seconds}{Cycle}$ $ET = IC \times CPI \times CT$ Frequency(i.e., clock rate)

 $1GHz = 10^9Hz = \frac{1}{10^9}sec \ per \ cycle = 1 \ ns \ per \ cycle$ 



## **Performance Equation (X)**

 Assume that we have an application composed with a total of 500000000 instructions, in which 20% of them are "Type-A" instructions with an average CPI of 8 cycles, 20% of them are "Type-B" instructions with an average CPI of 4 cycles and the rest instructions are "Type-C" instructions with average CPI of 1 cycle. If the processor runs at **3 GHz**, how long is the execution time?

A. 3.67 sec	$ET = (5 \times 10^9) \times (20\% \times 8 + 20\% \times 4 + 6\%)$
B. 5 sec	average CPI
C. 6.67 sec	$ET = IC \times CPI \times CT$
D. 15 sec	

E. 45 sec



 $60\% \times 1) \times \frac{1}{3 \times 10^{-9}} sec = 5$ 



• The relative performance between two machines, X and Y. Y is n times faster than X

$$n = \frac{Execution \ Time_X}{Execution \ Time_Y}$$

• The speedup of Y over X

$$Speedup = \frac{Execution \ Time_X}{Execution \ Time_Y}$$

## **Speedup of Y over X**

 Consider the same program on the following two machines, X and Y. By how much Y is faster than X?

	Clock Rate	Instructions	Percentage of Type-A Insts.	CPI of Type-A Insts.	Percentage of Type-B Insts.	CPI of Type-B Insts.	Percentage of Type-C Insts.	CPI of Type-C Insts.
Machine X	3 GHz	500000	20%	8	20%	4	60%	1
Machine Y	5 GHz	500000	20%	13	20%	4	60%	1
A.	0.2 <i>ET</i>	$T_Y = (5 \times 1)$	$0^{6}) \times (20\%)$	$\times 13 + 20^{\circ}$	$\% \times 4 + 60$	$\% \times 1) \times \frac{1}{5}$	$\frac{1}{\times 10^{-9}}sec =$	4
B.	0.25 <i>Sp</i>	$eedup = \frac{E}{F}$	Execution Time <sub>X</sub> Execution Time <sub>x</sub>	-				
<b>C</b> .	0.8	5	$\frac{5}{2}$ - 1.05					
D.	1.25	4	- = 1.23					
E.	No chang	ges						



## What Affects Each Factor in Performance Equation

## **Summary of CPU Performance Equation**



- IC (Instruction Count)
  - ISA, Compiler, algorithm, programming language, programmer
- CPI (Cycles Per Instruction)
  - Machine Implementation, microarchitecture, compiler, application, algorithm, programming language, programmer
- Cycle Time (Seconds Per Cycle)
  - Process Technology, microarchitecture, programmer

## **Programmer's impact**

• By adding the "sort" in the following code snippet, what the programmer changes in the performance equation to achieve **better** performance? std::sort(data, data + arraySize);

```
for (unsigned c = 0; c < arraySize*1000; ++c) {</pre>
        if (data[c%arraySize] >= INT_MAX/2)
             sum ++;
    }
```

- B. IC we increased IC, suppose to make the
- C. CT performance worse

D. IC & CPI

A. CP



## Demo — programmer & performance

for(i = 0; i < ARRAY\_SIZE; i++)</pre>  $\{$ for(j = 0; j < ARRAY\_SIZE; j++)</pre> c[i][j] = a[i][j]+b[i][j]; } }

for(j = 0; :
{
 for(i = 0
 {
 c[i][j]
 }
 }
}

$O(n^2)$	Complexity		
Same	Instruction Count?		
Same	<b>Clock Rate</b>		
Better	CPI		

- for(j = 0; j < ARRAY\_SIZE; j++)</pre>
  - for(i = 0; i < ARRAY\_SIZE; i++)</pre>

c[i][j] = a[i][j]+b[i][j];

 $O(n^2)$ 







### **Programmers can also set the cycle time**

https://software.intel.com/sites/default/files/comment/1716807/how-to-change-frequency-on-linux-pub.txt

```
_____
Subject: setting CPU speed on running linux system
```

If the OS is Linux, you can manually control the CPU speed by reading and writing some virtual files in the "/proc"

```
1.) Is the system capable of software CPU speed control?
If the "directory" /sys/devices/system/cpu/cpu0/cpufreq exists, speed is controllable.
-- If it does not exist, you may need to go to the BIOS and turn on EIST and any other C and P state control and vi:
```

```
2.) What speed is the box set to now?
Do the following:
$ cd /sys/devices/system/cpu
$ cat ./cpu0/cpufreg/cpuinfo max freq
3193000
$ cat ./cpu0/cpufreg/cpuinfo_min_freg
1596000
3.) What speeds can I set to?
Do
$ cat /sys/devices/system/cpu/cpu0/cpufreg/scaling available frequencies
It will list highest settable to lowest; example from my NHM "Smackover" DX58SO HEDT board, I see:
3193000 3192000 3059000 2926000 2793000 2660000 2527000 2394000 2261000 2128000 1995000 1862000 1729000 159600
You can choose from among those numbers to set the "high water" mark and "low water" mark for speed. If you set "h:
4.) Show me how to set all to highest settable speed!
Use the following little sh/ksh/bash script:
$ cd /sys/devices/system/cpu # a virtual directory made visible by device drivers
$ newSpeedTop=`awk '{print $1}' ./cpu0/cpufreg/scaling available frequencies`
$ newSpeedLow=$newSpeedTop # make them the same in this example
$ for c in ./cpu[0-9]* ; do
   echo $newSpeedTop >${c}/cpufreg/scaling max freq
>
   echo $newSpeedLow >${c}/cpufreg/scaling min freg
>
> done
ŝ
5.) How do I return to the default - i.e. allow machine to vary from highest to lowest?
Edit line # 3 of the script above, and re-run it. Change the line:
$ newSpeedLow=$newSpeedTop # make them the same in this example
```



## **Programming languages**

How many instructions are there in "Hello, world!"

	Instruction count	LOC	Ranking
С	600k	6	1
C++	ЗМ	6	2
Java	~210M	8	5
Perl	10M	4	3
Python	~30M	1	4





### **Source Code**

### Program

0f00bb27 509cbd23 00005d24 Data 0000bd24 2ca422a0 130020e4 00003d24 2ca4e2b3

00c2e800 80000008 00c2f000 80000008 00c2f800 00000008 00c30000 00000008

Compiler

(e.g., gcc)

### Recap: How my "Java code" becomes a "program"



# Jave Bytecode (.class)

Compiler

**Source Code** 

Ë

lava

00c2e800 cafebabe 00000008 0000033 001d0a00 00c2f000 ta 06000f09 80000008 00100011 00c2f800 0800120a 00000008 00130014 00c30000 07001507 80000008 ne Time Cost!

### **Recap: How my "Python code" becomes a "program"**



## **Source Code** python Perl

### Interpreter (e.g., python)

### Program

0f00bb27 509cbd23 00005d24 0000bd24

00c2e800 80000008 00c2f000 80000008

00c30000

00000008

### **Revisited the demo with compiler optimizations!**

- gcc has different optimization levels.
  - -O0 no optimizations
  - -O3 typically the best-performing optimization

```
for(i = 0; i < ARRAY_SIZE; i++)</pre>
     \mathbf{I}
       for(j = 0; j < ARRAY_SIZE; j++)</pre>
c[i][j] = a[i][j]+b[i][j];
```





### for(j = 0; j < ARRAY\_SIZE; j++)</pre> for(i = 0; i < ARRAY\_SIZE; i++)</pre> c[i][j] = a[i][j]+b[i][j];

### **Demo revisited — compiler optimization**

- Compiler can reduce the instruction count, change CPI — with "limited scope"
- Compiler CANNOT help improving "crummy" source code

if(option) std::sort(data, data + arraySize); **Compiler can never add this — only the programmer can!** for (unsigned c = 0; c < arraySize\*1000; ++c) {</pre> if (data[c%arraySize] >= INT MAX/2) sum ++;

### How many operations: CISC v.s. RISC

- CISC (Complex Instruction Set Computing)
  - Examples: x86, Motorola 68K
  - Provide many powerful/complex instructions
    - Many: more than 1503 instructions since 2016
    - Powerful/complex: an instruction can perform both ALU and memory operations
    - Each instruction takes more cycles to execute
- RISC (Reduced Instruction Set Computer)
  - Examples: ARMv8, RISC-V, MIPS (the first RISC instruction, invented by the authors of our textbook)
  - Each instruction only performs simple tasks
  - Easy to decode
  - Each instruction takes less cycles to execute



### **RISC-V v.s. x86**

- Using the same language, the same source code, regarding the compiled program on x86 and RISC-V, how many of the following statements is/are "generally" correct?
  - ① The RISC-V version would contain more instructions than its x86 version
  - <sup>2</sup> The RISC-V version tends to incur fewer memory accesses than its x86 version
  - ③ The RISC-V version needs a processor with higher clock rate than its x86 version if the CPI of both versions are similar
  - ④ The RISC-V version needs a processor with lower CPI than its x86 version if the x86 processor runs at the same clock rate
  - A. 0
  - **B**. 1
  - C. 2
  - D. 3

## Amdahl's Law — and It's Implication in the Multicore Era

H&P Chapter 1.9 M. D. Hill and M. R. Marty. Amdahl's Law in the Multicore Era. In Computer, vol. 41, no. 7, pp. 33–38, July 2008.

### **Amdahl's Law**



 $Speedup_{enhanced}(f, s) = \frac{1}{(1-f) + \frac{f}{s}}$ 

f — The fraction of time in the original program s — The speedup we can achieve on f









enhanced

Execution Time<sub>enhanced</sub> =  $(1-f) + f/s \leftarrow$ 

$$Speedup_{enhanced} = \frac{Execution Time_{baseline}}{Execution Time_{enhanced}}$$

$$\frac{1}{f) + \frac{f}{s}}$$



 $\frac{c_{baseline}}{c_{enhanced}} = \frac{1}{(1-f) + \frac{f}{s}}$ 

## **Practicing Amdahl's Law**

- Final Fantasy XV spends lots of time loading a map within which period that 95% of the time on the accessing the H.D.D., the rest in the operating system, file system and the I/O protocol. If we replace the H.D.D. with a flash drive, which provides 100x faster access time. By how much can we speed up the map loading process?
  - A. ~7x





I run this game from an 7200 RPM hardrive and load times are pretty long... do anyone run this game form an SSD? are load times good?

### **Amdahl's Law on Multiple Optimizations**

- We can apply Amdahl's law for multiple optimizations •
- These optimizations must be dis-joint!
  - If optimization #1 and optimization #2 are dis-joint:



$$\frac{1}{t_1 - f_{Opt2}} + \frac{f\_Opt1}{s\_Opt1} + \frac{f\_Opt2}{s\_Opt2}$$

### Practicing Amdahl's Law (2)

 Final Fantasy XV spends lots of time loading a map — within which period that 95% of the time on the accessing the H.D.D., the rest in the operating system, file system and the I/O protocol. If we replace the H.D.D. with a flash drive, which provides 100x faster access time and a better processor to accelerate the software overhead by 2x. By how much can we speed up the map loading process?



- A. ~7x
- B. ~10x



 $Speedup_{enhanced}(95\%, 5\%, 100, 2) = \frac{1}{(1-95\%)}$ 

E. ~100x

27

$$\frac{1}{-5\%) + \frac{95\%}{100} + \frac{5\%}{2}} = 28.98 \times$$

### Amdahl's Law Corollary #1

The maximum speedup is bounded by

$$Speedup_{max}(f, \infty) = \frac{1}{(1-f) + \frac{f}{\infty}}$$
$$Speedup_{max}(f, \infty) = \frac{1}{(1-f)}$$



## **Speedup further!**

 With the latest flash memory technologies, the system spends 16% of time on accessing the flash, and the software overhead is now 84%. If we want to adopt a new memory technology to replace flash to achieve 2x speedup on loading maps, how much faster the new technology needs to be?



- A. ~5x
- B. ~10x
- C. ~20x

 $Speedup_{max}(16\%,\infty) = \frac{1}{(1-16\%)} = 1.19$ 

D. ~100x

E. None of the above

### 2x is not possible

### **Corollary #1 on Multiple Optimizations**

If we can pick just one thing to work on/optimize •

$Speedup_{max}(f_1, \infty)$ :	$=\frac{1}{(1-c)}$
$C = 1 \qquad (C)$	$(1 - f_1)$ 1
$Speedup_{max}(f_2, \infty)$ :	$=$ $(1 - f_2)$
Speedup $(f_2, \infty)$ :	
$C = 1 \qquad (C)$	$(1 - f_3)$
$Speedup_{max}(f_4, \infty)$ :	$=\frac{1}{(1-f_4)}$



### **1-f<sub>1</sub>-f<sub>2</sub>-f<sub>3</sub>-f<sub>4</sub>**

### The biggest $f_x$ would lead to the largest *Speedup<sub>max</sub>*!

### **Corollary #2 — make the common case fast!**

- When f is small, optimizations will have little effect.
- Common == most time consuming not necessarily the most frequent
- The uncommon case doesn't make much difference
- The common case can change based on inputs, compiler options, optimizations you've applied, etc.

fect. essarily the most

erence ts, compiler

### If we repeatedly optimizing our design based on Amdahl's law...

### **Storage Media**



### CPU

- With optimization, the common becomes uncommon.
- An uncommon case will (hopefully) become the new common case.
- Now you have a new target for optimization.
- You have to revisit "Amdahl's Law" every time you applied some optimization

Moneta: A High-Performance Storage Array Architecture for Next-Generation, Non-volatile Memories Adrian M. Caulfield, Arup De, Joel Coburn, Todor I. Mollov, Rajesh K. Gupta, and Steven Swanson Proceedings of the 2010 43rd Annual IEEE/ACM International Symposium on Microarchitecture, 2010.



### **Don't hurt non-common part too mach**

- If the program spend 90% in A, 10% in B. Assume that an optimization can accelerate A by 9x, by hurts B by 10x...
- Assume the original execution time is T. The new execution time  $ET_{new} = \frac{ET_{old} \times 90\%}{0} + ET_{old} \times 10\% \times 10$  $ET_{new} = 1.1 \times ET_{old}$  $Speedup = \frac{ET_{old}}{ET_{even}} = \frac{ET_{old}}{1.1 \times ET_{even}} = 0.91 \times \dots \text{slowdown!}$

You may not use Amdahl's Law for this case as Amdahl's Law does NOT (1) consider overhead (2) bound to slowdown

### Amdahl's Law on Multicore Architectures

• Symmetric multicore processor with *n* cores (if we assume the processor performance scales perfectly)

$$Speedup_{parallel}(f_{parallelizable}, n) = \frac{1}{(1 - f_{parallel})}$$



## Corollary #3, Corollary #4 & Corollary #5

$$Speedup_{parallel}(f_{parallelizable}, \infty) = \frac{1}{(1 - f_{parallel})}$$
$$Speedup_{parallel}(f_{parallelizable}, \infty) = \frac{1}{(1 - f_{parallel})}$$

- Single-core performance still matters it will eventually dominate the performance
- Finding more "parallelizable" parts is also important
- If we can build a processor with unlimited parallelism the complexity doesn't matter as long as the algorithm can utilize all parallelism that's why bitonic sort works!







## "Fair" Comparisons

Andrew Davison. Twelve Ways to Fool the Masses When Giving Performance Results on Parallel Computers. In Humour the Computer, MITP, 1995


### Is TFLOPS (Tera FLoating-point Operations Per Second) a good metric?

 $TFLOPS = \frac{\# of floating point instructions \times 10^{-12}}{Exection Time}$ 

 $IC \times \%$  of floating point instructions  $\times 10^{-12}$ 

 $IC \times CPI \times CT$ 

% of floating point instructions  $\times 10^{-12}$ 

 $\overline{CPI \times CT}$ 

- Cannot compare different ISA/compiler
  - What if the compiler can generate code with fewer instructions?
  - What if new architecture has more IC but also lower CPI?
- Does not make sense if the application is not floating point intensive

### IC is gone!



### 12 ways to Fool the Masses When Giving Performance Results on Parallel Computers

- Quote only 32-bit performance results, not 64-bit results.
- Present performance figures for an inner kernel, and then represent these figures as the performance of the entire application.
- Quietly employ assembly code and other low-level language constructs.
- Scale up the problem size with the number of processors, but omit any mention of this fact.
- Quote performance results projected to a full system.
- Compare your results against scalar, unoptimized code on Crays.
- When direct run time comparisons are required, compare with an old code on an obsolete system.
- If MFLOPS rates must be quoted, base the operation count on the parallel implementation, not on the best sequential implementation.
- Quote performance in terms of processor utilization, parallel speedups or MFLOPS per dollar.
- Mutilate the algorithm used in the parallel implementation to match the architecture.
- Measure parallel run times on a dedicated system, but measure conventional run times in a busy environment.
- If all else fails, show pretty pictures and animated videos, and don't talk about performance.

- e on an obsolete system. I implementation, not on
- <sup>.</sup> MFLOPS per dollar. rchitecture. onal run times in a busy

	=		nvidia.com		Ċ			
Artificial Intelligence Computing Leadership from N CLOUD & DATA CENTER	PRODUCTS -	SOLUTIONS 🔻	APPS	▼ FOR DE\	/ELOPERS			
Tesla V100				AI TRAINING	AI INFERENCE			

### Deep Learning Training in Less Than a Workday



Server Config: Dual Xeon E5-2699 v4 2.6 GHz | 8X NVIDIA® Tesla® P100 or V100 | ResNet-50 Training on MXNet for 90 Epochs with 1.28M ImageNet Dataset.

### **AI TRAINING**

From recognizing speech to training virtual personal assistants and teaching autonomous cars to drive, data scientists are taking on increasingly complex challenges with AI. Solving these kinds of problems requires training deep learning models that are exponentially growing in complexity, in a practical amount of time.

With 640 Tensor Cores, Tesla V100 is the world's first GPU to break the 100 teraFLOPS (TFLOPS) barrier of deep learning performance. The next generation of NVIDIA NVLink<sup>™</sup> connects multiple V100 GPUs at up to 300 GB/s to create the world's most powerful computing servers. AI models that would consume weeks of computing resources on previous systems can now be trained in a few days. With this dramatic reduction in training time, a whole new world of problems will now be solvable with AI.



### TECHNOLOGIES -

### E HPC DATA CENTER GPUS SPECIFICATIONS

### The Most Advanced Data Center GPU Ever Built.

NVIDIA® Tesla® V100 is the world's most advanced data center. GPU ever built to accelerate AI, HPC, and graphics. Powered by NVIDIA Volta, the latest GPU architecture, Tesla V100 offers the performance of up to 100 CPUs in a single GPU—enabling data scientists, researchers, and engineers to tackle challenges that were once thought impossible.





### 1 GPU Node Replaces Up To 54 CPU Nodes Node Replacement: HPC Mixed Workload

Max Power

### SPECIFICATIONS



Tesla V100 PCle



### Tesla V100 SXM2

GPU Architecture	NVIDIA	Volta
NVIDIA Tensor Cores	64	10
NVIDIA CUDA <sup>«</sup> Cores	5,1	20
Double-Precision Performance	7 TFLOPS	7.8 TFLOPS
Single-Precision Performance	14 TFLOPS	15.7 TFLOPS
Tensor Performance	112 TFLOPS	125 TFLOPS
GPU Memory	32GB /16	GB HBM2
Memory		_ ,

### 900GB/sec

	Ye	5
	32GB/sec	300GB/sec
face	PCIe Gen3	NVIDIA NVLink
	PCIe Full Height/Length	SXM2

# They try to tell it's the better AI hardware

https://blogs.nvidia.com/blog/2017/04/10/ai-drives-rise-accelerated-computing-datacenter/

	K80 2012	TPU 2015	P40 2016
Inferences/Sec <10ms latency	1/ <sub>13</sub> X	1X	2X
Training TOPS	6 FP32	NA	12 FP32
Inference TOPS	6 FP32	90 INT8	48 INT8
On-chip Memory	16 MB	24 MB	11 MB
Power	300W	75W	250W
Bandwidth	320 GB/S	34 GB/S	350 GB/S



# What's wrong with inferences per second?

- There is no standard on how they inference
  - What model?
  - What dataset?
- That's why Facebook is trying to promote an AI benchmark **MLPerf** ٠

is an inaccurate summary performance metric. Our results show that IPS is a poor overall performance summary for NN hardware, as it's simply the inverse of the complexity of the typical inference in the application (e.g., the number, size, and type of NN layers). For example, the TPU runs the 4-layer MLP1 at 360,000 IPS but the 89-layer CNN1 at only 4,700 IPS, so TPU IPS vary by 75X! Thus, using IPS as the single-speed summary is even more misleading for NN accelerators than MIPS or FLOPS are for regular processors [23], so IPS should be even more disparaged. To compare NN machines better, we need a benchmark suite written at a high-level to port it to the wide variety of NN architectures. Fathom is a promising new attempt at such a benchmark suite [3].

Pitfall: For NN hardware, Inferences Per Second (IPS)



Node Replacement: HPC Mixed Workload

1 GPU Node Replaces Up To 54 CPU Nodes

CPU Server: Dual Xeon Gold 6140@2.30GHz, GPU Servers: same CPU server w/ 4x V100 PCIe | CUDA Version: CUDA 9.x| Dataset: NAMD (STMV), GTC (mpi#proc.in), MILC (APEX Medium), SPECFEM3D [four\_material\_simple\_model] | To arrive at CPU node equivalence, we use measured benchmark with up to 8 CPU nodes. Then we use linear scaling to scale beyond 8 nodes.

### HIGH PERFORMANCE COMPUTING (HPC)

HPC is a fundamental pillar of modern science. From predicting weather to discovering drugs to finding new energy sources, researchers use large computing systems to simulate and predict our world. AI extends traditional HPC by allowing researchers to analyze large volumes of data for rapid insights where simulation alone cannot fully predict the real world.

Tesla V100 is engineered for the convergence of AI and HPC. It offers a platform for HPC systems to excel at both computational science for scientific simulation and data science for finding insights in data. By pairing NVIDIA CUDA<sup>®</sup> cores and Tensor Cores within a unified architecture, a single server with Tesla V100 GPUs can replace hundreds of commodity CPU-only servers for both traditional HPC and AI workloads. Every researcher and engineer can now afford an AI supercomputer to tackle their most challenging work.

# Choose the right metric — Latency v.s. Throughput/Bandwidth

# Latency v.s. Bandwidth/Throughput

- Latency the amount of time to finish an operation
  - access time
  - response time
- Throughput the amount of work can be done within a given period of time
  - bandwidth (MB/Sec, GB/Sec, Mbps, Gbps)
  - IOPs
  - MFLOPs



## The performance between RAID and SSD

 Compare (X) RAID consists of 4x H.D.D. where each has 10 ms access time and 125 MB/sec bandwidth — aggregated bandwidth at 500 MB/Sec (Y) a single SSD with 100 us access time and 550MB/Sec bandwidth. Both accept 4KB data as the smallest request size. If we want to load a program with 100KB code size, how much faster is Y over X at least?

B. 1.1x  $ET_{HDD_{BestCase}} = 10 ms$ C. 4x  $ET_{SSD_{warst}} = \frac{100KB}{4K} \times 100 \ us = 2.5 \ ms$ D. 4.4x E. 100x

## Latency and Bandwidth trade-off

- Increase bandwidth can hurt the response time of a single task
- If you want to transfer a 2 Peta-Byte video from UCSD
  - 100 miles (161 km) from UCR
  - Assume that you have a 100Gbps ethernet
    - 2 Peta-byte over 167772 seconds = 1.94 Days
    - 22.5TB in 30 minutes
    - Bandwidth: 100 Gbps



# Memory Gap Problem



### **Moore's Law**<sup>(1)</sup>

### The establishment Reliability coun degree of integration will be achieved with linear Present and future Interneted electronics is antablished In almost e ICs are widely applicable By integrated electronics, I mean Increasing the yield demonstrated h to integrated electronics in the linear are or nev technologies which are referred to There is no fundamental obstacle to achieving and level of production-low compared to that of distronics today as well as any additidevice yields of 100%. At present, packaging costs vable result in electronics functions suppli crete components-it offers reduced systems cost, so far exceed the cost of the semiconductor structure itself that there is no incentive to improve illo, fo ICs are increasingly p and in many systems improved performance has yields, but they can be raised as high as is ecohe reli been realized. ICs are more reliable nomically justified. No barrier exists comparable that to miniaturize electronics equipment to the thermodynamic equilibrium considerations ilure as the creasingly complex electronic functi Heat problem space with minimum weight. Sever Will it be possible to remove the heat generated taoish by tens of thousands of components in a single evolved, including microassembly individual components, thin-film silicon chipi 15 eat is a solvable issue OF THE COMPONENTS TED FUNCTION semiconductor integrated circuits. Moore' 13 12 Day of reckoning Impo Clearly, we will be able to build such component-10 Two-mil squares crammed equipment. Next, we ask under what 9 circumstances we should do it. The total cost of historic With the dimensional tolerances already being making a particular system function must be miniemployed in integrated circuits, isolated high-per **ICs are easy to manufacture** mized. To do so, we could amortize the engineerformance transistors can be built on centers two ing over several identical items, or evolve flexible thousandths of an inch apart. Such a two-mil square techniques for the engineering of large functions and they're getting smaller and can also contain several kilohms of resistance so that no disproportionate expense need be borne ICs are small by a particular array. Perhaps newly devised desmaller esigning ICs can be easy proponents onto integrated circuits', Electronics 38 (8). (1) Mo YEAR 50

### Linear circuitry

Integration will not change linear systems as radically as digital systems. Still, a considerable

### **Performance gap between Processor/Memory**





## Performance of modern DRAM

			Best case ac	cess time (no pr	echarge)	Precharge needed
Production year	Chip size	DRAM type	RAS time (ns)	CAS time (ns)	Total (ns)	Total (ns)
2000	256M bit	DDR1	21	21	42	63
2002	512M bit	DDR1	15	15	30	45
2004	1G bit	DDR2	15	15	30	45
2006	2G bit	DDR2	10	10	20	30
2010	4G bit	DDR3	13	13	26	39
2016	8G bit	DDR4	13	13	26	39

Figure 2.4 Capacity and access times for DDR SDRAMs by year of production. Access time is for a random memory word and assumes a new row must be opened. If the row is in a different bank, we assume the bank is precharged; if the row is not open, then a precharge is required, and the access time is longer. As the number of banks has increased, the ability to hide the precharge time has also increased. DDR4 SDRAMs were initially expected in 2014, but did not begin production until early 2016.



# The impact of "slow" memory

- Assume that we have a processor running @ 2 GHz and a program with 30% of load/store instructions. If the computer has "perfect" memory, the CPI is just 1. Now, consider we have DDR4 and the program is wellbehaved that precharge is never necessary — the access latency is simply 26 ns. What's the average CPI (pick the most close one)?
  - A. 9
  - B. 17
  - C. 27
  - D. 35

E. 69

$$1 + 100\% \times (52) + 30\% \times 52$$

### Fetch Instructio





### 68.6 cycles

### **Alternatives?**

Memory technology	Typical access time	
SRAM semiconductor memory	0.5–2.5ns	
DRAM semiconductor memory	50–70ns	
Flash semiconductor memory	5,000-50,000ns	
Magnetic disk	5,000,000-20,000,000ns	
	<b>Fast, but expens</b>	ive

### **\$ per GiB in 2012** \$500-\$1000 \$10-\$20 \$0.75-\$1.00

### \$0.05-\$0.10





### How can memory hierarchy help in performance?

 Assume that we have a processor running @ 2 GHz and a program with 30% of load/store instructions. If the computer has "perfect" memory, the CPI is just 1. Now, in addition to DDR4, whose latency 26 ns, we also got an SRAM cache with latency of just at 0.5ns and can capture 90% of the desired data/instructions. what's the average CPI (pick the most close one)?

A. 2  
B. 4  
C. 8 
$$1 + (1 - 90\%) \times [100\% \times (52) + 30\%$$
  
D. 16  
E. 32

### $\times$ 52] = 7.76 cycles cess Data

# L1? L2? L3?



CPU Cache	s Mainboa	ard Memo	ory S
Processor			
Name		Intel Core	e i7 970
Code Name	Coffe	e Lake	Max
Package		Socket 1	151 LG
Technology	14 nm	Core V	oltage
Specification	In	tel® Core	™ i7-97
Family	6	Mo	odel
Ext. Family	6	Ext. Mo	odel
Instructions	MMX, SSE, AES, AVX,	SSE2, SSE AVX2, FMA	E3, SSS \3, TSX
Clocks (Core	#0)		-Cach
Core Speed	4798.85	5 MHz	L1D
Multiplier	x 48.0 (8	3 - 49 )	L1 In
Bus Speed	99.98	MHz	Leve
Rated FSB			Leve
Selection	Socket #1	<b>v</b>	





### How can deeper memory hierarchy help in performance?

- Assume that we have a processor running @ 2 GHz and a program with 30% of load/store instructions. If the computer has "perfect" memory, the CPI is just 1. Now, in addition to DDR4, whose latency 26 ns, we also got a 2-level SRAM caches with
  - it's 1st-level one at latency of 0.5ns and can capture 90% of the desired data/ instructions.
  - the 2nd-level at latency of 5ns and can capture 60% of the desired data/instructions What's the average CPI (pick the most close one)?

A. 2  $1 + (1 - 90\%) \times [10 + (1 - 60\%) \times 52 + 30\% \times (10 + (1 - 60\%) \times 52)] = 5 \ cycle$ B. 4 C. 8 **Fetch Instructions Access Data** D. 16 E. 32

# Why adding small SRAMs would work?



- Spatial locality application tends to visit nearby stuffs in the memory
  - Code the current instruction, and then PC + 4
  - Data the current element in an array, then the next
- Temporal locality application revisit the same thing again and again
  - Code loops, frequently invoked functions
  - Data the same data can be read/write many times

# Most of time, your program is just visiting a very small amount of data/instructions within a given window

# Architecting the Cache

### Processor Load/store only access a "word" each time

Load 0x000A

Core

Registers

0x0000																								—
0x1000	AAAA	BBBB	CCCC	DDDD	EEEE	FFFF	GGGG	НННН	AAAA	BBBB	CCCC	DDDD	EEEE	FFFF	GGGG	НННН	AAAA	BBBB	CCCC	DDDD	EEEE	FFFF	GGGG	н
0x2000	AAAA	BBBB	CCCC	DDDD	EEEE	FFFF	GGGG	нннн	ΑΑΑΑ	BBBB	сссс	DDDD	EEEE	FFFF	GGGG	нннн	АААА	BBBB	CCCC	DDDD	EEEE	FFFF	GGGG	н
0x3000																								
0x4000																								
0x5000																								
0x6000																								
0x7000																								
0x8000																								
	•	•	•	•	• •	•	•	•	•	• •	• •	•	•	•	•	• •	• •	• •	• •	• •	• •	• •	• •	•
	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	
	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	• •	
	• •	• •	• •	• •	• •	• •	• • •	• • •	• •	• •	• • •	• •	• • •	• •	• •	• •	• • •	• •	• •	• • •	• •	• •	• • •	•
	• • •	• •	• • •	• •	• •	• • •	• • •	• •	• • •	• •	• •	• •	• •	• •	• • •	• • •	• •	• •	• •	• • •	• •	• • •	• • •	
	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	



egi	ste	rs		oad 0x000A SRABB CCD												"Lo me	ogio emo	call ory "bl	lly" partition y space into blocks"												
AAAA	BBBB	CCCC	DDDD	EEEE	FFFF	GGGG	нннн	AABB	CCDD	EEFF	GGHH	EEEE	FFFF	GGGG	НННН	AAAA	BBBB	CCCC	DDDD	EEEE	FFFF	GGGG	ннн	АААА	BBBB	CCCC	DDDD	EEE	FFFF	GGGG	нн
AAAA	BBBB	CCCC	DDDD	EEEE	FFFF	GGGG	нннн	АААА	BBBB	CCCC	DDDD	EEEE	FFFF	GGGG	нннн	AAAA	BBBB	CCCC	DDDD	EEEE	FFFF	GGGG	нннн	AAAA	BBBB		טטטט	EEEE	FFFF	GGGG	HHł
:	•	:	:	:		•	•	:	:	:	:	•	•	:	:	:	:	:	:	•	•	:	•	:	•	:	:	:	:	:	:
•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
		•	•	•	•	•	•	•	•	•	•	•								•			•	•	•	•	•	•	•	•	



# How to tell who is there?



××××××××××××××××× 23456789ABCDEF AABBCCDDEEGGFFHH IIJJKKLLMMNNOOPP **IIJJKKLLMMNNOOPP** IIJJKKLLMMNNOOPP IIJJKKLLMMNNOOPP QQRRSSTTUUVVWWXX QQRRSSTTUUVVWWXX



# Tell if the block here can be used



### data 0123456789ABCDEF FFHH BBCCDDEE **IIJJKKLLMMNNOOPP QQRRSSTTUUVV**WWXX **YYZZAABBCCDDEEFF AABBCCDDEEGG**FFHH **IIJJKKLLMMNNOOPP QQRRSSTTUUVV**WWXX **YYZZAABBCCDDEEFF IIJJKKLLMMNNOOPP** QQRRSSTTUUVVWWXX **YYZZAABBCCDDEEFF AABBCCDDEEGG**FFHH **IIJJKKLLMMNNOOPP** QQRRSSTTUUVVWWXX **QQRRSSTTUUVV**WWXX **YYZZAABBCCDDEEFF**

### **Way-associative cache**





### C = ABS

- C: Capacity in data arrays
- A: Way-Associativity how many blocks within a set
  - N-way: N blocks in a set, A = N
  - 1 for direct-mapped cache
- **B**: **B**lock Size (Cacheline)
  - How many bytes in a block
- S: Number of Sets:
  - A set contains blocks sharing the same index
  - 1 for fully associate cache



# **Corollary of C = ABS**

set block index offset tag 0b000100000100100

memory address:

- number of bits in block offset lg(B)
- number of bits in set index: lg(S)
- tag bits: address\_length lg(S) lg(B)
  - address\_length is 32 bits for 32-bit machine
- (address / block size) % S = set index



# **AMD Phenom II**

- L1 data (D-L1) cache configuration of AMD Phenom II
  - Size 64KB, 2-way set associativity, 64B block
  - Assume 64-bit memory address
  - Which of the following is correct?
    - A. Tag is 49 bits
    - B. Index is 8 bits
    - C. Offset is 7 bits
    - D. The cache has 1024 sets
    - E. None of the above

- 64KB = 2 \* 64 \* S
  - S = 512
- offset = lg(64) = 6 bits
- index = lg(512) = 9 bits
- tag = 64 lg(512) lg(64) = 49 bits

# intel Core i7

- L1 data (D-L1) cache configuration of Core i7
  - Size 32KB, 8-way set associativity, 64B block
  - Assume 64-bit memory address
  - Which of the following is NOT correct?
    - A. Tag is 52 bits
    - B. Index is 6 bits
    - C. Offset is 6 bits
    - D. The cache has 128 sets

C = ABS

32KB = 8 \* 64 \* S

S = 64

offset = lg(64) = 6 bits

index = lg(64) = 6 bits

tag = 64 - lg(64) - lg(64) = 52 bits
# Put everything all together: How cache interacts with CPU





## **Performance evaluation considering cache**

- If the load/store instruction hits in L1 cache where the hit time is usually the same as a CPU cycle
  - The CPI of this instruction is the base CPI
- If the load/store instruction misses in L1, we need to access L2
  - The CPI of this instruction needs to include the cycles of accessing L2
- If the load/store instruction misses in both L1 and L2, we need to go to lower memory hierarchy (L3 or DRAM)
  - The CPI of this instruction needs to include the cycles of accessing L2, L3, DRAM

## How to evaluate cache performance

CPI<sub>Average</sub>: the average CPI of a memory instruction

CPI<sub>base</sub> + miss\_rate<sub>L1</sub>\*miss\_penalty<sub>L1</sub> CPI<sub>Average</sub>=

miss\_penalty<sub>L1</sub>=  $CPI_{accessing L2}$ +miss\_rate<sub>L2</sub>\*miss\_penalty<sub>L2</sub>

miss\_penalty<sub>L2</sub> =  $CPI_{accessing_L3} + miss_rate_{L3} * miss_penalty_{L3}$ 

miss\_penalty<sub>L3</sub> = CPl<sub>accessing DRAM</sub>+miss\_rate<sub>DRAM</sub>\*miss\_penalty<sub>DRAM</sub>

 If the problem is asking for average memory access time, transform the CPI values into/from time by multiplying with CPU cycle time!



# **Cache & Performance**







=

# **Cause of cache misses**



# **3Cs of misses**

- Compulsory miss
  - Cold start miss. First-time access to a block
- Capacity miss
  - The working set size of an application is bigger than cache size
- Conflict miss
  - Required data replaced by block(s) mapping to the same set
  - Similar collision in hash

# Simulate the cache!



# **Tips for cache simulation**

- Figure out the memory access patterns
  - Address sequences from your code
  - The behavior/locality of the variables/arrays
- Partition the address
  - Use C=ABS
  - Find out tag, index
- Check your current cache content
  - Hit: for the same index, if you can find the same tag there.
  - Otherwise, miss
    - Compulsory misses: you never accessed the same (tag, index) pair before
    - Conflict misses: the tag appeared in the same index before
    - Replace the least recently used block with the requesting block



# Simulate a direct-mapped cache

- Consider a direct mapped (1-way) cache with 256 bytes total capacity, a block size of 16 bytes, and the application repeatedly reading the following memory addresses:
  - Ob100000000, Ob100001000, Ob1000010000, Ob1000010100, 0b1100010000
    - C = A B S
    - S=256/(16\*1) = 16
    - lg(16) = 4 : 4 bits are used for the index
    - lg(16) = 4 : 4 bits are used for the byte offset
    - The tag is 48 (4 + 4) = 40 bits
    - For example: 0b1000 0000 0000 0000 0000 0000 1000 0000

tag







# Simulate a direct-mapped cache

	V	D	Tag	Data
0	1	0	0b10	
1	1	0	0b10	
2	0	0		
3	0	0		
4	0	0		
5	0	0		
6	0	0		
7	0	0		
8	0	0		
9	0	0		
10	0	0		
11	0	0		
12	0	0		
13	0	0		
14	0	0		
15	0	0		

tag index

0b1	-0	0000	000
0b1	-0	0000	100
0b1	-0	0001	000
0b1	-0	0001	010
0b1	1	0001	000
0b1	-0	0000	000
0b1	-0	0000	100
0b1	-0	0001	000
0b1	.0	0001	010



**compulsory miss** 90 hit! 90 **compulsory miss** 90 90 hit! 90 **compulsory miss** 90 hit! 90 hit! 90 conflict miss 90 hit!

# Simulate a 2-way cache

- Consider a 2-way cache with 256 bytes total capacity, a block size of 16 bytes, and the application repeatedly reading the following memory addresses:
  - Ob100000000, Ob100001000, Ob1000010000, 0b1000010100, 0b1100010000
    - C = A B S
    - S=256/(16\*2) = 8
    - $8 = 2^3$ : 3 bits are used for the index
    - $16 = 2^4 : 4$  bits are used for the byte offset
    - The tag is 32 (3 + 4) = 25 bits
    - For example: 0b1000 0000 0000 0000 0000 0000 0001 0000

## tag





# Simulate a 2-way cache

	V	D	Tag	Data	V	D	Tag	Data
0	1	0	0b10		0	0		
1	1	0	0b10		1	0	0b11	
2	0	0			0	0		
3	0	0			0	0		
4	0	0			0	0		
5	0	0			0	0		
6	0	0			0	0		
7	0	0			0	0		

tag index 0b10 0000 0b10 0001 0b10 0000 0b10 0000 0b10 0001 0b10 0001



## 0b10 0000 0000 compulsory miss hit! 1000 0b10 0001 0000 compulsory miss 0100 hit! 0b11 0001 0000 compulsory miss 0000 hit! 1000 hit! 0000 hit 0100 hit!

- D-L1 Cache configuration of AMD Phenom II
  - Size 64KB, 2-way set associativity, 64B block, LRU policy, write-allocate, write-back, and assuming 32-bit address.

```
int a[16384], b[16384], c[16384];
/* c = 0x10000, a = 0x20000, b = 0x30000 */
for(i = 0; i < 512; i++) {
    c[i] = a[i] + b[i];
   //load a, b, and then store to c
}
```

What's the data cache miss rate for this code?

A. 6.25% C = ABSB. 56.25% 64KB = 2 \* 64 \* S C. 66.67% S = 512offset = lg(64) = 6 bits 68.75% D index = lg(512) = 9 bits E. 100% tag = 64 - lg(512) - lg(64) = 49 bits

Size 64KB, 2-way set associativity, 64B block, LRU policy, write-allocate, write-back, and assuming 48-bit address.  $\bullet$ 

int /\* for

a[16384], c = 0x1000 (i = 0; i	C = ABS 64KB = 2 * 64 * S S = 512 offset = lg(64) = 6 bits				
c[i] = a[	index = lg(512) = 9 bits tag = the rest bits				
	address in hex	address in binary tag index offset	tag	index	hit? miss?
load a[0]	0x20000	0b10 0000 0000 00 <mark>00 0000</mark>	0x4	0	compulsory miss
load b[0]	0x30000	0b11 0 <mark>000 0000 00</mark> 00 0000	0x6	0	compulsory miss
store c[0]	0x10000	<mark>0b01 0</mark> 000 0000 00 <mark>00 0000</mark>	0x2	0	compulsory miss, evict
load a[1]	0x20004	0b10 0000 0000 00 <mark>00 0100</mark>	0x4	0	conflict miss, evict 0x6
load b[1]	0x30004	<mark>0b11 0</mark> 000 0000 00 <mark>00 0100</mark>	0x6	0	conflict miss, evict 0x2
store c[1]	0x10004	<mark>0b01 0</mark> 000 0000 00 <mark>00 0100</mark>	0x2	0	conflict miss, evict 0x4
load a[15]	0x2003C	<mark>0b10 0</mark> 000 0000 00 <mark>11 1100</mark>	0x4	0	miss, evict 0x6
load b[15]	0x3003C	<mark>0b11 0</mark> 000 0000 00 <mark>11 1100</mark>	0x6	0	miss, evict 0x2
store c[15]	0x1003C	<mark>0b01 0</mark> 000 0000 00 <mark>11 1100</mark>	0x2	0	miss, evict 0x4
load a[16]	0x20040	0b10 0000 0000 0100 0000	0x4	1	compulsory miss
load b[16]	0x30040	0b11 0000 0000 01 <mark>00 0000</mark>	0x6	1	compulsory miss
store c[16]	0x10040	0b01 0000 0000 0100 0000	0x2	1	compulsory miss, evict
		89			

## 100% miss rate!

- D-L1 Cache configuration of AMD Phenom II
  - Size 64KB, 2-way set associativity, 64B block, LRU policy, write-allocate, write-back, and assuming 32-bit address.

```
int a[16384], b[16384], c[16384];
/* c = 0x10000, a = 0x20000, b = 0x30000 */
for(i = 0; i < 512; i++) {
    c[i] = a[i] + b[i];
   //load a, b, and then store to c
}
```

What's the data cache miss rate for this code?

A. 6.25% C = ABSB. 56.25% 64KB = 2 \* 64 \* S C. 66.67% S = 512offset = lg(64) = 6 bits 68.75% D. index = lg(512) = 9 bits E. 100% tag = 64 - lg(512) - lg(64) = 49 bits

# Matrix transpose

```
double A[16384], B[16384];
int N=128;
for(i = 0; i < N; i++)
    for(j = 0; j < N; j++)
        B[i*N+j] = A[j*N+i];
// assume load A[j*N+i] and then store B[i*N+j]
// &A[0] is 0x20000, &B[0] is 0x40000
```

## What's the access sequence of A[] looks like?

A[0], A[128], A[256], ..., A[127\*128], A[1], A[129]..., A[127\*128+1], ...

What's the access sequence of B[] looks like? B[0], B[1], B[2], .....

# Improving 3Cs

# Improvement of 3Cs

- 3Cs and A, B, C of caches
  - Compulsory miss
    - Increase B: increase miss penalty (more data must be fetched from lower hierarchy)
  - Capacity miss
    - Increase C: increase cost, access time, power
  - Conflict miss
    - Increase A: increase access time and power
- Or modify the memory access pattern of your program!



# Programming and memory performance

# Memory addressing/alignment

- Almost every popular ISA architecture uses "byte-addressing" to access memory locations
- Instructions generally work faster when the given memory address is aligned
  - Aligned if an instruction accesses an object of size n at address X, the access is aligned if  $X \mod n = 0$ .
  - Some architecture/processor does not support aligned access at all
  - Therefore, compilers only allocate objects on "aligned" address



# The result of sizeof(struct student)

• Consider the following data structure: struct student { int id; double \*homework; int participation; double midterm; double average; };



64-bit

What's the output of

printf("%lu\n",sizeof(struct student))?

- A. 20
- B. 28
- C. 32
- D. 36





# Loop interchange/fission/fusion

# Demo — programmer & performance

for(i = 0; i < ARRAY\_SIZE; i++)</pre>  $\{$ for(j = 0; j < ARRAY\_SIZE; j++)</pre> c[i][j] = a[i][j]+b[i][j]; } }

for(j = 0;
{
 for(i = 0
 {
 c[i][j]
 }
}

 $O(n^2)$ ComplexitySameInstruction Count?SameClock RateBetterCPI

- for(j = 0; j < ARRAY\_SIZE; j++)</pre>
  - for(i = 0; i < ARRAY\_SIZE; i++)</pre>

c[i][j] = a[i][j]+b[i][j];

 $O(n^2)$ 







- D-L1 Cache configuration of AMD Phenom II
  - Size 64KB, 2-way set associativity, 64B block, LRU policy, write-allocate, write-back, and assuming 32-bit address.

```
int a[16384], b[16384], c[16384];
/* c = 0x10000, a = 0x20000, b = 0x30000 */
for(i = 0; i < 512; i++) {
    c[i] = a[i] + b[i];
    //load a, b, and then store to c
}
```

What's the data cache miss rate for this code?

A. 6.25% C = ABSB. 56.25% 64KB = 2 \* 64 \* S C. 66.67% S = 512offset = lg(64) = 6 bits 68.75% D. index = lg(512) = 9 bits 100% tag = 64 - lg(512) - lg(64) = 49 bits

# What if the code look like this?

- D-L1 Cache configuration of AMD Phenom II
  - Size 64KB, 2-way set associativity, 64B block, LRU policy, write-allocate, write-back, and assuming 32-bit address.

```
int a[16384], b[16384], c[16384];
/* c = 0x10000, a = 0x20000, b = 0x30000 */
for(i = 0; i < 512; i++)
    c[i] = a[i]; //load a and then store to c
for(i = 0; i < 512; i++)
    c[i] += b[i]; //load b, load c, add, and then store to c
```

What's the data cache miss rate for this code?

- A. 6.25%
- B. 56.25%
- C. 66.67%
- 68.75% D.





# **Loop Fusion**

```
/* Before */
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
        a[i][j] = 1/b[i][j] * c[i][j];
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
        d[i][j] = a[i][j] + c[i][j];
/* After */
for (i = 0; i < N; i = i+1)
    for (j = 0; j < N; j = j+1)
    {
        a[i][j] = 1/b[i][j] * c[i][j];
         d[i][j] = a[i][j] + c[i][j];
    }
```

2 misses per access to a & c vs. one miss per access

# Blocking

# **Case study: Matrix Multiplication**

```
for(i = 0; i < ARRAY_SIZE; i++) {</pre>
                                         Algorithm class tells you it's O(n<sup>3</sup>)
 for(j = 0; j < ARRAY_SIZE; j++) {</pre>
   for(k = 0; k < ARRAY_SIZE; k++) {</pre>
                                           If n=512, it takes about 1 sec
     c[i][j] += a[i][k]*b[k][j];
   }
 }
                                   How long is it take when n=1024?
}
```



# **Matrix Multiplication**



- If each dimension of your matrix is 1024
  - Each row takes 1024\*8 bytes = 8KB
  - The L1 \$ of intel Core i7 is 32KB, 8-way, 64-byte blocked
  - You can only hold at most 4 rows/columns of each matrix!
  - You need the same row when j increase!

## Very likely a miss if array is large





## **Block algorithm for matrix multiplication**

- Discover the cache miss rate
  - valgrind --tool=cachegrind cmd
    - cachegrind is a tool profiling the cache performance
  - Performance counter
    - Intel<sup>®</sup> Performance Counter Monitor http://www.intel.com/software/pcm/



## **Block algorithm for matrix multiplication**



## You only need to hold these sub-matrices in your cache

for(kk = k; kk < k+(ARRAY\_SIZE/n); kk++)</pre>



# **Matrix Transpose**

```
// Transpose matrix b into b_t
                                                                        b_t[i][j] += b[j][i];
                                                                    }
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {</pre>
                                                                  }
  for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {</pre>
    for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {</pre>
         for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)</pre>
           for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)</pre>
             for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)</pre>
               c[ii][jj] += a[ii][kk]*b[kk][jj];
                                                                       }
                                                                    }
```

for(i = 0; i < ARRAY\_SIZE; i+=(ARRAY\_SIZE/n)) {</pre> for(j = 0; j < ARRAY\_SIZE; j+=(ARRAY\_SIZE/n)) {</pre>

for(i = 0; i < ARRAY\_SIZE; i+=(ARRAY\_SIZE/n)) {</pre> for(j = 0; j < ARRAY\_SIZE; j+=(ARRAY\_SIZE/n)) {</pre> for(k = 0; k < ARRAY\_SIZE; k+=(ARRAY\_SIZE/n)) {</pre> for(ii = i; ii < i+(ARRAY\_SIZE/n); ii++)</pre> for(jj = j; jj < j+(ARRAY\_SIZE/n); jj++)</pre> for(kk = k; kk < k+(ARRAY\_SIZE/n); kk++)</pre> // Compute on b\_t c[ii][jj] += a[ii][kk]\*b\_t[jj][kk];

# Prefetching
### **Characteristic of memory accesses**

```
for(i = 0;i < 1000000; i++) {</pre>
     D[i] = rand();
}
```





### Prefetching

```
for(i = 0;i < 1000000; i++) {</pre>
     D[i] = rand();
     // prefetch D[i+8] if i % 8 == 0
}
```



# Prefetching

- Identify the access pattern and proactively fetch data/ instruction before the application asks for the data/instruction
  - Trigger the cache miss earlier to eliminate the miss when the application needs the data/instruction
- Hardware prefetch
  - The processor can keep track the distance between misses. If there is a pattern, fetch miss\_data\_address+distance for a miss
- Software prefetch
  - Load data into XO
  - Using prefetch instructions

### Demo

- x86 provide prefetch instructions
- As a programmer, you may insert mm prefetch in x86 programs to perform software prefetch for your code
- gcc also has a flag "-fprefetch-loop-arrays" to automatically insert software prefetch instructions

### Where can prefetch work effectively?

 How many of the following code snippet can "prefetching" effectively help improving performance?





# Advanced Hardware Techniques in Improving Memory Performance

### Without banks



### **Multibanks & non-blocking caches**





### return block **0**xDEAEBE



### **Pipelined access and multi-banked caches**

 Assume each bank in the \$ takes 10 ns to serve a request, and the \$ can take the next request 1 ns after assigning a request to a bank — if we have 4 banks and we want to serve 4 requests, what's the speedup over non-banked, non-pipelined \$? — pick the closest one



E. 5x

 $ET_{baseline} = 4 \times 10 \ ns = 40 \ ns$  $ET_{banked} = 10 \ ns + 3 \times 1 \ ns = 13 \ ns$  $Speedup = \frac{Execution Time_{baseline}}{Execution Time_{banked}}$  $=\frac{40}{13}=3.08\times$ 

### **Early Restart and Critical Word First**

- Don't wait for full block to be loaded before restarting CPU Early restart—As soon as the requested word of the block arrives, send it to the CPU and let the CPU continue execution
  - Critical Word First—Request the missed word first from memory and send it to the CPU as soon as it arrives; let the CPU continue execution while filling the rest of the words in the block. Also called wrapped fetch and requested word first
- Most useful with large blocks
- Spatial locality a problem; often we want the next sequential word soon, so not always a benefit (early restart).

# **Midterm Logistics**



### **For midterm**

- No cheat sheet allowed
- No cheating allowed
- We will have some problems require you to write
- You may bring a calculator
- You should bring pen/pencil/eraser
- My last office hour before midterm this Wednesday @ 1pm-2pm friday is cancelled.

### Format of the midterm

- Multiple choices \* 10 like your clicker/reading quizzes multiple choices questions
- Short answer question \* 5
  - Each answer MUST be less than 30 words
  - Writing more than 30 words is equivalent to writing 0 words
- Homework style free-answer questions \* 3
  - You need to clearly write down the original form of the applied equation/formula
  - You need to replace each term accordingly with numbers
  - You will have some credits for right equations even though the final number isn't correct
  - You will receive 0 credits if we only see the numbers



# Sample Midterm

### Identify the performance bottleneck

 Why does an Intel Core i7 @ 3.5 GHz usually perform better than an Intel Core i5 @ 3.5 GHz or AMD FX-8350@4GHz?



- A. Because the instruction count of the program are different
- B. Because the clock rate of AMD FX is higher
- C. Because the CPI of Core i7 is better
- D. Because the clock rate of AMD FX is higher and CPI of Core i7 is better
- E. None of the above



Sysbench 2014 from http://www.anandtech.com/

### Amdahl's Law on Multicore Architectures

- Regarding Amdahl's Law on multicore architectures, how many of the following statements is/are correct?
  - If we have unlimited parallelism, the performance of each parallel piece does not matter as long as the performance slowdown in each piece is bounded
  - ② With unlimited amount of parallel hardware units, single-core performance does not matter anymore
  - ③ With unlimited amount of parallel hardware units, the maximum speedup will be bounded by the fraction of parallel parts
  - ④ With unlimited amount of parallel hardware units, the effect of scheduling and data exchange overhead is minor
  - A. 0
  - B. 1
  - C. 2
  - D. 3
  - E. 4

### How programmer affects performance?

- Performance equation consists of the following three factors
  - ① IC
  - $\bigcirc$  CPI
  - 3 CT

How many can a **programmer** affect?

- A. 0
- **B**. 1
- C. 2
- D. 3



### **Demo** — programmer & performance

```
for(i = 0; i < ARRAY_SIZE; i++)</pre>
      for(j = 0; j < ARRAY_SIZE; j++)</pre>
c[i][j] = a[i][j]+b[i][j];
    }
```

- How many of the following make(s) the performance of A better than **B**?
  - ① IC
  - $\bigcirc$  CPI
  - 3 CT
  - A. 0
  - **B**. 1
  - C. 2
  - D. 3

- RRAY\_SIZE; j++)
- ARRAY\_SIZE; i++)
- i][j]+b[i][j];

## **Fair comparison**

- How many of the following comparisons are fair?
  - ① Comparing the frame rates of Halo 5 on AMD RyZen 1600X and civilization on Intel Core i7 7700X
  - ② Using bit torrent to compare the network throughput on two machines
  - Comparing the frame rates of Halo 5 using medium settings on AMD RyZen 3 1600X and low settings on Intel Core i7 7700X
  - ④ Using the peak floating point performance to judge the gaming performance of machines using AMD RyZen 1600X and Intel Core i7 7700X
  - A. 0
  - **B**. 1
  - C. 2
  - D. 3
  - E. 4



 Which description about locality of arrays sum and A in the following code is the most accurate? for(i = 0; i < 100000; i++){

```
sum[i\%10] += A[i];
```

}

- A. Access of A has temporal locality, sum has spatial locality
- B. Both A and sum have temporal locality, and sum also has spatial locality
- C. Access of A has spatial locality, sum has temporal locality
- D. Both A and sum have spatial locality
- E. Both A and sum have spatial locality, and sum also has temporal locality

# **3Cs and A, B, C**

- Regarding 3Cs: compulsory, conflict and capacity misses and A, B, C: associativity, block size, capacity How many of the following are correct?
  - ① Increasing associativity can reduce conflict misses
  - ② Increasing associativity can reduce hit time
  - Increasing block size can increase the miss penalty 3
  - Increasing block size can reduce compulsory misses (4)
  - A. 0
  - **B**. 1
  - C. 2
  - D. 3

### E. 4

## intel Core i7

- L1 data (D-L1) cache configuration of Core i7
  - Size 32KB, 8-way set associativity, 64B block
  - Assume 64-bit memory address
  - Which of the following is NOT correct?
    - A. Tag is 52 bits
    - B. Index is 6 bits
    - C. Offset is 6 bits
    - D. The cache has 128 sets

### Virtual indexed, physical tagged cache limits the cache size

- If you want to build a virtual indexed, physical tagged cache with 32KB capacity, which of the following configuration is possible? Assume the system use 4K pages.
  - A. 32B blocks, 2-way
  - B. 32B blocks, 4-way
  - C. 64B blocks, 4-way
  - D. 64B blocks, 8-way



### When we have virtual memory...

- In a modern x86-64 processor supports virtual memory through, how many memory accesses can an instruction incur?
  - A. 2
  - B. 4
  - C. 6
  - D. 8
  - E. More than 10



### Sample short answer questions (< 30 words)

- What is RISC? What is CISC? List two pros/cons for each
- What are the limitations of compiler optimizations? Can you list two?
- Please define Amdahl's Law and explain each term in it
- Please define the CPU performance equation and explain each term.
- Can you list two things affecting each term in the performance equation?
- What's the difference between latency and throughput? When should you use latency or throughput to judge performance?
- What's "benchmark" suite? Why is it important?
- Why TFLOPS or inferences per second is not a good metrics?

### **Performance Equation/Speedup**

• Assume that we have an application composed with a total of 500000 instructions, in which 20% of them are the load/store instructions with an average CPI of 6 cycles, and the rest instructions are integer instructions with average CPI of 1 cycle. If the processor runs at 1GHz, how long is the execution time? If hardware technology improves the processor clock rate to 2GHz, but making load/store CPI to be 12 cycles, how much is the speedup?

### Amdahl's Law for multiple optimizations

- Assume that memory access takes 30% of execution time.
  - Cache can speedup 80% of memory operation by a factor of 4
  - L2 cache can speedup 50% of the remaining 20% by a factor of 2
- What's the total speedup?

### ecution time. by a factor of 4 % by a factor of 2

### Performance evaluation with cache

• Consider the following cache configuration on RISC-V processor:

	I-L1	D-L1	L2
size	32K	32K	256K
block size	64 Bytes	64 Bytes	64 Bytes
associativity	2-way	2-way	8-way
access time	1 cycle (no penalty if it's a hit)	1 cycle (no penalty if it's a hit)	10 cycles
local miss rate	2%	10%, 20% dirty	15% (i.e., 15% of L1 mis also miss in the L2), 30
Write policy	N/A	Write-back, write allocate	
Replacement	LRU replacement policy		

The application has 20% branches, 10% loads/stores, 70% integer instructions. Assume that TLB miss rate is 2% and it requires 100 cycles to handle a TLB miss. Also assume that the branch predictor has a hit rate of 87.5%, what's the CPI of branch, L/S, and integer instructions? What is the average CPI?



### DRAM

Big enough 4KB pages

100 cycles

sses, % dirty

### **Cache simulation**

• The processor has a 8KB, 256B blocked, 2-way L1 cache. Consider the following code:

```
for(i=0;i<256;i++) {</pre>
     a[i] = b[i] + c[i];
// load a[i] and load b[i], store to c[i]
// \&a[0] = 0 \times 10000, \&b[0] = 0 \times 20000, \&c[0] = 0 \times 30000
}
```

- What's the total miss rate? How many of the misses are compulsory misses? How many of the misses are conflict misses?
- How can you improve the cache performance of the above code through changing hardware?
- How can you improve the performance **without** changing hardware?