

# CS 110

# Computer Architecture

## *Advanced Caches*

Instructor:

Sören Schwertfeger and Chundong Wang

<https://robotics.shanghaitech.edu.cn/courses/ca/20s>

School of Information Science and Technology SIST

ShanghaiTech University

Slides based on UC Berkeley's CS61C and  
Carnegie Mellon Univ. ECE447 (2015)

# Midterm

- Date: Tuesday, May 26th, 2020
- Time: 10:15- 12:15 (normal lecture slot++)
  - Be there latest 10:00 – **we start 10:15 sharp!**
- Venue: 4 rooms – check on egate which room you are!
  - SPST1-503
  - SPST1-201
  - SPST1-501
  - SIST1A-106
- Closed book:
  - You can bring **two** A4 pages with notes (both sides; in **English**): Write your Chinese and **Pinyin** name on the top! **Handwritten** by you!
    - **Final**: you can bring **three** A4 pages
  - You will be provided with the RISC-V "green sheet"
  - No other material allowed!

# Midterm I

- Wear your Corona mask! =>
- Switch cell phones **off!**  
(not silent mode – off!)
  - Put them in your bags.
- Bags under the table. Nothing except paper, pen, 1 drink, 1 snack, your student ID card on the table!
- No other electronic devices are allowed!
  - No ear plugs, music, smartwatch...
- Anybody touching any electronic device will **FAIL** the course!
- Anybody found cheating (copy your neighbors answers, additional material, ...) will **FAIL** the course!



Admin



Admin



Admin



Admin

# COMPUTER ORGANIZATION AND DESIGN

THE HARDWARE/SOFTWARE INTERFACE

 RISC-V EDITION



**MK**  
MORGAN KAUFMANN

DAVID A. PATTERSON  
JOHN L. HENNESSY

FUNCTIONS OF SEVERAL VARIABLES		$z = f(x, y)$ $w = f(x, y, z)$	DOMAIN: Allowed $(x, y)$ , $(x, y, z)$ RANGES: $z, w$ 's
<b>LEVEL CURVES</b> Contour Maps (2-D) Surface Layers (3-D)	<b>FUNCTION OF n VARIABLES</b> $\mathbb{R}^n \rightarrow \mathbb{R}$ 1. As a function of $n$ real variables $x_1, x_2, \dots, x_n$ 2. As a function of a single variable $x$ , $x_1, \dots, x_n$ 3. As a function of a single vector var. $\vec{x} = (x_1, \dots, x_n)$	<b><math>\mathbb{R}^3</math> DEFINITION OF CONTINUITY</b> LET $f$ BE A FUNCTION OF 2 VARIABLES DEFINED ON A DISK W/ CENTER $(a, b)$ , EXCEPT POSSIBLY $(a, b)$ . THEN $\lim_{(x, y) \rightarrow (a, b)} f(x, y) = L$ IF FOR EVERY $\epsilon > 0$ , THERE IS A CORRESPONDING $\delta > 0$ ST. IF $(x, y) \in D_f$ AND $\sqrt{(x-a)^2 + (y-b)^2} < \delta$	
<b>PARTIAL DERIVATIVES</b>	Derivatives w/ respect to one variable, while holding the other variables constant SAME METHODS FOR FUNCTIONS OF MORE THAN TWO VARIABLES	IF THE LIMIT AS A FUNCTION APPROACHES A POINT $(a, b)$ ALONG TWO DIFFERENT PATHS IS NOT THE SAME, THE LIMIT DOES NOT EXIST $\exists$ $f(x, y)$ IS CONTINUOUS AT $(a, b)$ IF THE LIMIT OF $(x, y)$ AS $(x, y) \rightarrow (a, b)$ EXISTS.	
<b>SECOND PARTIAL DERIVATIVES</b>	<b>CLAIRAUT'S THEOREM</b> IF $f_{xy}$ AND $f_{yx}$ ARE BOTH CONTINUOUS $f_{xy}(a, b) = f_{yx}(a, b)$ <b>PARTIAL DIFF. EQS</b> LAPLACE'S EQUATION $\Delta u = 0$ etc... HEAT EQUATION $\Delta u = c$ etc... WAVE EQUATION $\Delta u = c \frac{\partial^2 u}{\partial t^2}$	COMPOSITE FUNCTIONS OF CONTINUOUS FUNCTIONS ARE CONTINUOUS, AS ARE SUMS AND PRODUCTS <b>EQUATIONS OF TANGENT PLANES TO SURFACES</b> $z = f(x, y)$ @ $(x_0, y_0, z_0)$ EVALUATED AT A POINT $z - z_0 = f_x(x_0, y_0)(x - x_0) + f_y(x_0, y_0)(y - y_0)$ <b>TOTAL DIFFERENTIAL</b> ( $dy = f'(x) dx$ SINGLE VARIABLE) $dz = f_x(x, y) dx + f_y(x, y) dy = \frac{\partial z}{\partial x} dx + \frac{\partial z}{\partial y} dy$ INCREMENTS $\Delta x, \Delta y, \Delta z$ DIFFERENTIALS $dx, dy, dz$ FOR SMALL $\Delta x, \Delta y$ $\Delta z \approx dz$ ( $\Delta z$ CHANGE IN HEIGHT OF SURFACE ( $\Delta z$ ) CHANGE IN HEIGHT OF THE TANGENT PLANE ( $dz$ ))	
<b>THE CHAIN RULE</b>	<b>SINGLE VARIABLE</b> $y = f(x)$ , $w = g(t)$ , IF $z = f(g(t))$ $\frac{dz}{dt} = \frac{df}{dg} \cdot \frac{dg}{dt}$ <b>CASE 1</b> $z = f(x, y)$ , $x = g(t)$ , $y = h(t)$ $\Rightarrow z = f(g(t), h(t))$ $\frac{dz}{dt} = \frac{\partial z}{\partial x} \frac{dx}{dt} + \frac{\partial z}{\partial y} \frac{dy}{dt}$ OR $w = f(x, y, z)$ , $x = g(t)$ , $y = h(t)$ , $z = k(t)$ $\Rightarrow w = f(g(t), h(t), k(t))$ $\frac{dw}{dt} = \frac{\partial w}{\partial x} \frac{dx}{dt} + \frac{\partial w}{\partial y} \frac{dy}{dt} + \frac{\partial w}{\partial z} \frac{dz}{dt}$ <b>CASE 2</b> $z = f(x, y)$ , $x = g(s, t)$ , $y = h(s, t)$ , $z = k(s, t)$ $\frac{dz}{ds} = \frac{\partial z}{\partial x} \frac{dx}{ds} + \frac{\partial z}{\partial y} \frac{dy}{ds} + \frac{\partial z}{\partial z} \frac{dz}{ds}$ <b>CHAIN RULE: GENERAL VERSION</b> $u = f(x_1, \dots, x_n)$ , $x_j = g_j(t_1, \dots, t_m)$ $\frac{du}{dt_i} = \frac{\partial u}{\partial x_1} \frac{dx_1}{dt_i} + \dots + \frac{\partial u}{\partial x_n} \frac{dx_n}{dt_i}$ FOR EACH $i = 1, 2, \dots, m$	IF $f_x$ AND $f_y$ ARE CONTINUOUS $\Delta z \approx dz$ ( $\Delta z$ CHANGE IN HEIGHT OF SURFACE ( $\Delta z$ ) CHANGE IN HEIGHT OF THE TANGENT PLANE ( $dz$ )) <b>THEOREM</b> $\Delta z = f_x(a, b) \Delta x + f_y(a, b) \Delta y + \epsilon$ $\epsilon$ AND $\epsilon_2$ ARE FUNCTIONS OF $\Delta x$ AND $\Delta y$ THAT APPROACH 0 AS $(\Delta x, \Delta y) \rightarrow (0, 0)$ DEF. $f$ IS DIFFERENTIABLE @ $(a, b)$ FOR $z = f(x, y)$ <b>DIFFERENTIAL DIAGRAM</b> (CASE 1 & 2) YOU CAN FIND DERIVATIVES FOR ALL THE TEMPERATURE SUBSTITUTIONS $\frac{dz}{dt_i} = \frac{\partial z}{\partial x} \frac{dx}{dt_i} + \frac{\partial z}{\partial y} \frac{dy}{dt_i} + \frac{\partial z}{\partial z} \frac{dz}{dt_i}$ $\frac{dz}{dt_i} = \frac{\partial z}{\partial x} \frac{dx}{dt_i} + \frac{\partial z}{\partial y} \frac{dy}{dt_i} + \frac{\partial z}{\partial z} \frac{dz}{dt_i}$	
<b>IMPLICIT DIFFERENTIATION</b>	You can always solve for $y$ and $dz$ $\frac{dz}{dx} = \frac{\partial z}{\partial x} - \frac{f_y}{f_x}$ $\frac{dz}{dy} = \frac{\partial z}{\partial y} - \frac{f_x}{f_y}$	<b>THE GRADIENT VECTOR</b> $\nabla f(x, y) = \left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$	
<b>TANGENT PLANE TO A LEVEL SURFACE</b> $\nabla f \cdot \vec{r} = 0$	$F_x(x - x_0) + F_y(y - y_0) + F_z(z - z_0) = 0$	<b>DIRECTIONAL DERIVATIVES</b> $\text{Diff } \vec{u} = (a, b)$ $\text{Diff } f(x, y) = f_x(x, y) a + f_y(x, y) b$ SAME FOR 3 VARIABLES $\text{Diff } f(x, y) = \nabla f(x, y) \cdot \vec{u}$ $\text{Diff } f$ MAX OCCURS WHEN $\vec{u}$ IS IN THE SAME DIR. AS $\nabla f$ $\text{Diff } f = \nabla f \cdot \vec{u} =  \nabla f   \vec{u}  \cos \theta =  \nabla f  \cos \theta$	
<b>NORMAL LINE TO LEVEL SURFACE</b> $\vec{r} = \frac{\nabla f}{ \nabla f }$	$\vec{r} = \frac{\nabla f}{ \nabla f }$	<b>THE GRADIENT VECTOR IS ORTHOGONAL TO THE LEVEL CURVES OF A SURFACE</b> $\nabla f$ HAS AS MANY COMPONENTS AS $f$ HAS INDEPENDENT VARIABLES. $N = (f_x, f_y, -1)$	
<b>SPECIAL CASE</b> $F_x(x, y) = 0$ , $F_y(x, y) = 0$ , $F_z(x, y) = 0$	LEVEL SURFACE W/ NO OLD DEFINITION THEN $F_x = -1$ , $F_y = (f_x, f_y, -1)$ AND TANGENT PLANE $z - z_0 = f_x(x - x_0) + f_y(y - y_0)$	<b>THE GRADIENT VECTOR POINTS IN THE DIRECTION OF STEEPEST ASCENT OR DESCENT (MAX/SURFACE)</b> <b>TO FIND THE NORMAL (AND LATER TANGENT) PLANE TO A SURFACE, LET THAT SURFACE BE THE LEVEL SET OF SOME HIGHER DIMENSIONAL FUNCTION. THEN THE GRADIENT OF THE HIGHER D. FUNCTION IS <math>\perp</math> TO YOUR SURFACE</b> $\nabla W = (2x, 2y, 2z)$ IS A NORMAL VECTOR TO THE $z = 2 - 3xy + z^2$	
<b>MAXIMUM AND MINIMUM VALUES</b> $z = f(x, y)$	$f_x(a, b) = 0$ , $f_y(a, b) = 0$ NECESSARY BUT NOT SUFFICIENT TO GUARANTEE A MAX. OR MIN. D = $\frac{f_{xx} f_{yy} - (f_{xy})^2}{(f_{xx})^2}$ $D > 0$ , $f_{xx} > 0$ LOCAL MIN. $D < 0$ , $f_{xx} < 0$ LOCAL MAX. $D = 0$ SADDLEPT. DO NOT USE.	<b>NEWTON'S METHOD</b> 1. A body's speed remains constant if the net force acting on it is zero (balanced). Acceleration is directly proportional to net force and inversely proportional to mass. Forces exist in action/reaction pairs. Remember that when gravity is acting against the sum of forces, it will be factored as a negative.	
<b>FINDING ABSOLUTE MAX. AND MINS. FOR A CLOSED BOUNDARY</b>	1. Find values of $f$ at the critical points of $f$ in $D$ . 2. Find the endpoint values of $f$ on the boundary of $D$ . 3. The largest value from 1, 2. IS THE ABS. MAX. THE SMALLEST IS THE ABS. MIN.	<b>CONVERSION</b> $m/s \times 3.6 = km/h$ $km/h \div 3.6 = m/s$	
<b>MAXIMIZING AND MINIMIZING</b> SET OF A FUNCTION OF TWO VARIABLES OF THE FORM $z = f(x, y)$ AND EACH DOMAIN USUAL CONSTRAINT			

**PHYSICS**

**WAVELENGTH**  $\lambda$   
 $v = \lambda f$   
SPEED OF LIGHT  $c = 3 \times 10^8$  m/s  
WAVELENGTH  $\lambda$   
SPEED OF WAVE  $v$   
FREQUENCY  $f$

**DIFFRACTION**  
DIFFRACTION IS THE BENDING OF LIGHT AROUND OBSTACLES AND THROUGH APERTURES.  
DIFFRACTION OCCURS WHEN THE SIZE OF THE APERTURE IS COMPARABLE TO THE WAVELENGTH.

**INTERFERENCE**  
CONSTRUCTIVE INTERFERENCE: when 2 or more waves travel in the same direction, the resulting wave is the sum of their displacements.  
DESTRUCTIVE INTERFERENCE: when 2 or more waves travel in opposite directions, the resulting wave is zero.

**REFLECTION**  
LAW OF REFLECTION: the angle of incidence is equal to the angle of reflection. The incident and reflected beams, and the normal, all lie on the same plane.

**REFRACTION**  
REFRACTION IS THE BENDING OF THE PATH OF LIGHT DUE TO A CHANGE IN SPEED AS IT ENTERS A MEDIUM OF DIFFERENT OPTICAL DENSITY.  
DENSE MEDIUM: high refractive index  
LESS DENSE MEDIUM: low refractive index  
SLOWER MEDIUM: refracted towards the normal  
FASTER MEDIUM: refracted away from the normal

**GRAVITY**  
GRAVITATIONAL POTENTIAL:  $U = mgh$   
GRAVITATIONAL POTENTIAL ENERGY:  $E_p = mgh$   
KINETIC ENERGY:  $E_k = \frac{1}{2}mv^2$   
POWER:  $P = \frac{W}{t}$

**VELOCITY**  $v$  (m/s)  
**DISPLACEMENT**  $s$  (m)  
**FORCE**  $F$  (N)  
**ACCELERATION**  $a$  (m/s<sup>2</sup>)  
**MASS**  $m$  (kg)  
**MOMENTUM**  $p$  (kg m/s)  
**IMPULSE**  $J$  (N s)  
**WORK**  $W$  (J)  
**ENERGY**  $E$  (J)  
**POWER**  $P$  (W)

**VECTOR ADDITION**  
Add vector quantities head to tail!

**GRAVITY**  
GRAVITY IS THE FORCE OF ATTRACTION BETWEEN TWO OBJECTS.  
GRAVITATIONAL FORCE:  $F_g = \frac{Gm_1m_2}{r^2}$   
GRAVITATIONAL ACCELERATION:  $g = 9.8$  m/s<sup>2</sup>

**NEWTON'S LAWS**  
1. A body's speed remains constant if the net force acting on it is zero (balanced).  
2. Acceleration is directly proportional to net force and inversely proportional to mass.  
3. Forces exist in action/reaction pairs.  
Remember that when gravity is acting against the sum of forces, it will be factored as a negative.

**OPTICS**  
OPACQUE: doesn't allow light to pass through.  
TRANSPARENT: allows light to pass through.  
TRANSLUCENT: allows light through but distorts its path (blur the image).  
NO MATERIAL CAN ALLOW 100% OF INCIDENT LIGHT TO PASS THROUGH.

**REFRACTION**  
REFRACTION IS THE BENDING OF THE PATH OF LIGHT DUE TO A CHANGE IN SPEED AS IT ENTERS A MEDIUM OF DIFFERENT OPTICAL DENSITY.  
DENSE MEDIUM: high refractive index  
LESS DENSE MEDIUM: low refractive index  
SLOWER MEDIUM: refracted towards the normal  
FASTER MEDIUM: refracted away from the normal

**THE SUN IS AN OURNING STAR**  
FUSION  
 $4H^1 \rightarrow He^4 + 2e^+ + 2\nu_e + \text{energy}$   
Energy is released when two light atoms fuse to form a heavier atom. This process releases energy when the strong nuclear force overcomes the repulsive force between the protons.  
The sun's main sequence energy is produced by the fusion of hydrogen into helium.

**STAR HEATS, SHINKS THEN EXPANDS**  
As the star heats up, it expands and its core contracts, increasing the temperature and pressure. This leads to the fusion of hydrogen into helium, which releases energy and causes the star to expand further.



### Old School Machine Structures

When C program starts

- C executable and is loaded into memory by OS (opening gates)
- OS put it on stack, then calls into application
- Run the top instruction every and the others
- Then will your processor read manual!

Valid Pointer Arithmetic

- Add an integer to a pointer
- Subtract 2 pointers for the same array
- Compare two pointers (C, C++, ...)
- Capture pointer to null
- Add two pointers / multiply

Five Fundamental Steps in Calling a Function

1. Put parameters in a place where function can access them
2. Transfer control to function
3. Access local storage
4. Execute needed part of the function
5. Put result value in a place where calling code can access it and restore any registers you used

### New School Machine Structures

Parallel Requests  
Parallel Threads  
Parallel Data  
Hardware descriptors

从高位 从下往上依次递增 总合名称

1. 位 左 → 右
2. 左 → 右
3. 左 → 右
4. 左 → 右
5. 左 → 右
6. 左 → 右
7. 左 → 右
8. 左 → 右
9. 左 → 右

int main (int argc, char \* argv[])

- argc contains the number of things on the command line
- argv is a pointer to an array containing the arguments as strings

Program's address space / Memory Address

Stack: local variables inside functions, grow downwards

Heap: space requested for dynamic data via malloc(), grows dynamically

Static data: variable declared outside functions, loaded when program starts, can be modified

Code: loaded when program starts, doesn't change

Return control to part of origin, since a function can be called from several parts in a program.

Save argument registers  
Save return address  
Save saved registers  
Local variables and functions

Stack  
↑  
Dynamic data  
↓  
Heap  
↓  
Reserved

### Guest Ideas in Computer Architecture

1. Abstraction
2. Moore's Law
3. Principle of Locality
4. Parallelism
5. Performance Measurement and Improvement
6. Dependability via Redundancy

Two's Complement Representation

- treats 0 as positive
- 32-bit word represents  $2^{32}$  into form  $-2^{31}$  to  $2^{31}-1$

### Components of a Computer

foo.c → cpp → fo.o → compiler

- cpp replaces comments with a double quote
- cpp comments begins with #

(standard type)

- char / signed char
- int
- long
- float
- double

10. 位级或 1 左 → 右
11. 位级或 1 左 → 右
12. 11 位级或
13. ? 右 → 左
14. = = = = 右 → 左
15. / 左 → 右

Stack

- free when function returns
- last in first out

Common Memory Problems

- Very uninitialized values
- Wrong memory that you did own
- Important use of free/malloc by memory with the pointer handle
- Memory leaks

Levels of Representation / Interpretation

- High level language
- Assembly language
- Machine language
- Machine / Interpretation
- Hardware Architecture Description
- Architecture / Implementation
- Logic Level Description

R format: used for instructions with immediate, lw and sw and branches (key and base)

Opcode rs rt rd shamt funct

opcode = 0

2 format: 5-bit field only represent number up to 31

If machine has multiple, then use at most 2 registers

used for lw, sw, and branches and with immediate

Label Upper lui

Imm \$to \$fs OADR ABCD

Opcode \$to \$fs OADR ABC

Opcode \$to \$fs OADR ABC

Opcode \$to \$fs OADR ABC



THE PURPOSE OF THIS REFERENCE GUIDE IS TO WALK THROUGH THE PROCESS OF BOOTING THE SIFT WORKSTATION, CREATING A TIMELINE ("SUPER" OR "MICRO") AND REVIEWING IT.

Admin

Computer-forensics11.sans.org/community/downloads  
 Download: SIFT Workstation VM Appliance  
 Download: SIFT Workstation Installation

2. BOOT SIFT VM  
 Login: sansforensics  
 Password: forensics

3. ELEVATE PRIVS  
 \$ sudo su

4. CONNECT IMAGE TO SIFT  
 Plug hard drive to physical host and attach to SIFT VM

- log2timeline PARSING PLUGINS
- apache2\_error - Apache2 error log file
  - chrome - Chrome history file
  - encase\_dirlisting - CSV file that is exported from encase
  - evt - Windows 2k/XP/2k3 Event Log
  - evtx - Windows Event Log File (EVTX)
  - exif - Metadata information from files using ExifTool
  - ff\_bookmark - Firefox bookmark file
  - firefox2 - Firefox 2 browser history
  - firefox3 - Firefox 3 history file
  - ftk\_dirlisting - CSV file that is exported from FTK Imager (dirlisting)
  - generic\_linux - Generic Linux logs that start with MMM.DD.HH:MM:SS
  - iehistory - index.dat file containing IE history
  - iis - IIS W3C log file
  - isatxt - ISA text export log file
  - jp\_ntfs\_change - CSV output file from JP (NTFS Change log)
  - mactime - Body file in the mactime format
  - mcafee - Log file
  - mft - NTFS MFT file
  - mssql\_errlog - ERRORLOG file produced by MS SQL server
  - ntuser - NTUSER.DAT registry file
  - opera - Opera's global history file
  - oxml - OpenXML document pcap
  - pcap - PCAP file
  - pdf - Available PDF document metadata
  - prefetch - Prefetch directory
  - recycler - Recycle bin directory
  - restore 0.9 - Restore point directory
  - safari - Safari History.plist file
  - sam - SAM registry file
  - security - SECURITY registry file
  - setupapi - SetupAPI log file in Windows XP
  - skype\_sql - Skype database
  - software - SOFTWARE registry file
  - sol - .sol (LSO) or a Flash cookie file
  - squid - Squid access log (http\_emulate off)
  - syslog - Linux Syslog log file
  - system - SYSTEM registry file
  - tlf - Body file in the TLF format
  - volatility - Volatility output files (psscan2, socksca2, ...)
  - win\_link - Windows shortcut file (or a link file)
  - wmiprov - wmiprov log file
  - xpfirewall - XP Firewall log

BY DAVID NIDES (12/16/2011)  
 TWITTER: @DAVNADS  
 BLOG: DAVNADS.BLOGSPOT.COM  
 EMAIL: DNIDES@KPMG.COM  
 CREDITS TO: ED GOINGS, ROB LEKSTON, KRISTINN GUDJONSSON, KPMG & PWC  
 QUESTIONS/FEEDBACK-CONTACT US

KEY  
 Red text – image/source  
 Blue text – mount point  
 Purple text – output file  
 Green text – log2timeline plugins  
 Brown text – TimeZone

5. HARD DRIVE MOUNTING (if you are using log2timeline-sift and Single DD you can skip to 7-A)

SINGLE OR SPLIT IMAGE (2 options):  
 # mount -t ntfs -o ro,loop,show\_sys\_files,streams\_interface=windows,offset=#### /mnt/ewf/<image> /mnt/windows\_mount/  
 # mount -t ntfs -o ro,loop,show\_sys\_files,streams\_interface=windows,offset=#### image.E01 /mnt/ewf/

MOUNT TO MOUNT POINT  
 # mount -t ntfs -o ro,loop,show\_sys\_files,streams\_interface=windows,offset=#### image.dd /mnt/windows\_mount/  
 # affuse image.001 /mnt/aff  
 # mount -t ntfs-3g -o loop,ro,show\_sys\_files,streams\_interface=windows,offset=#### /mnt/aff/<image> /mnt/windows\_mount/

HOW TO CALCULATE THE OFFSET FOR MOUNTING  
 1. Run mmls to query partition layout  
 # mmls image.E01  
 2. Identify partition and byte offset  
 3. (Partition byte offset) x (bytes per sector) = offset ##### to use!  
 Example: 63 X 512 = 32256

Note: If needed, repeat for each partition. Make new mount point:  
 # mkdir /mnt/windows\_mount2/

6. log2timeline default timezone is set to examiner local host. To change use -z [TIMEZONE] option. To list all available timezones:  
 # log2timeline -z list

7-A: AUTOMATED SUPER TIMELINE CREATION

log2timeline-sift -o -z [TIMEZONE] -p [PARTITION #] -i [IMAGE FILE]

DISK IMAGE (prompt for partition, mount, and run):  
 XP # log2timeline-sift -z EST5EDT -i image  
 WIN7 # log2timeline-sift -win7 -z EST5EDT -i image

FOR PARTITION (mount and run using all applicable plugins):  
 XP # log2timeline-sift -z EST5EDT -p 0 -i partition  
 WIN7 # log2timeline-sift -win7 -z EST5EDT -p 0 -i partition

OTHER USAGE EXAMPLES:  
 Display list of available plugins:  
 # log2timeline -f list  
 Run log2timeline using only specific plugins:  
 # log2timeline-sift -z EST5EDT -i image.dd  
 Help (man pages):  
 # log2timeline -h

8. CSV FILE OUTPUT (/cases/timeline-output-folder)

-date: Time of the event, in the format of MM/DD/YYYY  
 -time: Time of day, expressed in a 24h format, HH:MM:SS  
 -timezone: The timezone that was used to call the tool with.  
 -source: MACB meaning of the fields, comp w/ mactime format.  
 -sourcetype: Source short name (i.e. registry entries are REG)  
 -sourcetype: Desc of the source ("Internet Explorer" instead of WEBHIST)  
 -type: Timestamp type (i.e. "Last Accessed", "Last Written")  
 -user: Username associated with the entry, if one is available.  
 -host: Hostname associated with the entry, if one is available.  
 -short: Contains less text than the full description field.  
 -desc: where majority info is stored, the actual parsed desc of the entry.  
 -version: Version number of the timestamp object.  
 -filename: Filename with the full path that contained the entry  
 -inode: inode number of the file being parsed.  
 -notes: Some input modules insert additional information in the form of a note, which comes here. Or it can be used during the review.  
 -format: Input module name used to parse the file.  
 -extra: Additional information parsed is joined together and put here.

ANNUAL "MICRO" TIMELINE CREATION

log2timeline-sift -o -z [TIMEZONE] [-f FORMAT] [-z TIMEZONE] [-o OUTPUT MODULE] [-w FILE\_PATH] [-m METADATA] [-f FORMAT] [-z TIMEZONE] [-o OUTPUT MODULE] [-w FILE\_PATH] [-m METADATA]

EXTRACT METADATA (using log2timeline or fls)  
 # log2timeline -f mft -o mactime -r -z EST5EDT -w mft.body /mnt/volume/  
 OR Extract Metadata using Sleuthkit:  
 # fls -m "" -o offset -d > fls.body  
 Convert body file format to mactime:  
 # mactime -b fls.body -d > log2timeline.csv

ARTIFACTS (run I2I on mounted file system with plugins recursively)

Extract artifacts w/ log2timeline and run on mounted file system:  
 # log2timeline -f firefox3,chrome -o mactime -z EST5EDT -w web.body /mnt/volume/  
 Convert body file format to CSV format w/ mactime:  
 # mactime -b log2timeline.body -d > log2timeline.csv

9. FILTER TIMELINE

Filter timeline with date range to include only:  
 I2t\_process -b timeline.csv MM-DD-YYYY..MM-DD-YYYY > filtered.csv  
 Filter timeline with keyword list (one term per line in keywords.txt):  
 I2t\_process -b timeline.csv -k keywords.txt > filtered.csv  
 What sources are in your timeline?  
 awk -F, '{print \$6;}' timeline.csv | grep -v sourcetype | sort | uniq  
 Find all LNK files that reference E Drive  
 grep "Shortcut LNK" timeline.csv | grep "E:"  
 Find MountPoints2 entries that reference E Drive  
 grep "MountPoints2 key" timeline.csv | grep "E drive"  
 grep "USB timeline.csv" | grep "SetupAPIlog"

File System	M	A	C	B
Ext2/3	Modified	Accessed	Changed	N/A
FAT	Written	Accessed	N/A	Created
NTFS	File Modified	Accessed	MFT Modified	Created
UFS	Modified	Accessed	Changed	N/A

7-A & 7-B

HELP? OPTIONS? USAGE?  
 log2timeline -help  
 log2timeline-sift -help  
 I2t\_process -help

OTHER log2timeline OUTPUT FORMATS  
 Note: CSV is Default Output  
 -BeeDocs - Mac OS X visualization tool  
 -CEF - Common Event Format - ArcSight  
 -CFTL - XML file- CyberForensics TimeLab visualization tool  
 -CSV - comma separated value file  
 -Mactime - Both older and newer version of the format supported for use by TSK's mactime  
 -SIMILE - XML file - SIMILE timeline visualization widget  
 -SQLite - SQLite database  
 -TLN - Tab Delimited File  
 -TLN - Format used by some of H Carvey tools, expressed as a ASCII output  
 -TLNX - Format used by some of H Carvey tools, expressed as a XML document

10. CONNECT TO SIFT

- 1. SIFT Desktop -> SETTINGS -> OPTIONS -> Shared Folders -> Always Enabled (Check)
- 2. SIFT Desktop -> VMWare-Shared-Drive
- 3. Access from a Win Machine \\SIFTWORKSTATION

11. REVIEW TIMELINE

- Review timelines using:
- Open, Soft, Filter with Excel
  - Import into SPLUNK
  - SIMILE
  - Tapestry

# Content

- Main topics: Everything till (including) Lecture 16
  - Number representation (int & float)
  - C
  - CALL
  - RISC-V
  - SDS; Datapath & Control
  - Pipelining & Superscalar
  - Caches
- Plus general "Computer Architecture" knowledge

# Quiz on TLB

Piazza: “Video Lecture 22 VM”

- TLB Reach: Size of largest virtual address space that can be simultaneously mapped by TLB
- If we increase the page size, say, from 4KB to 16KB, can we increase the TLB reach?
  - A. Yes, of course
  - B. No, TLB reach is fixed given a processor and memory space

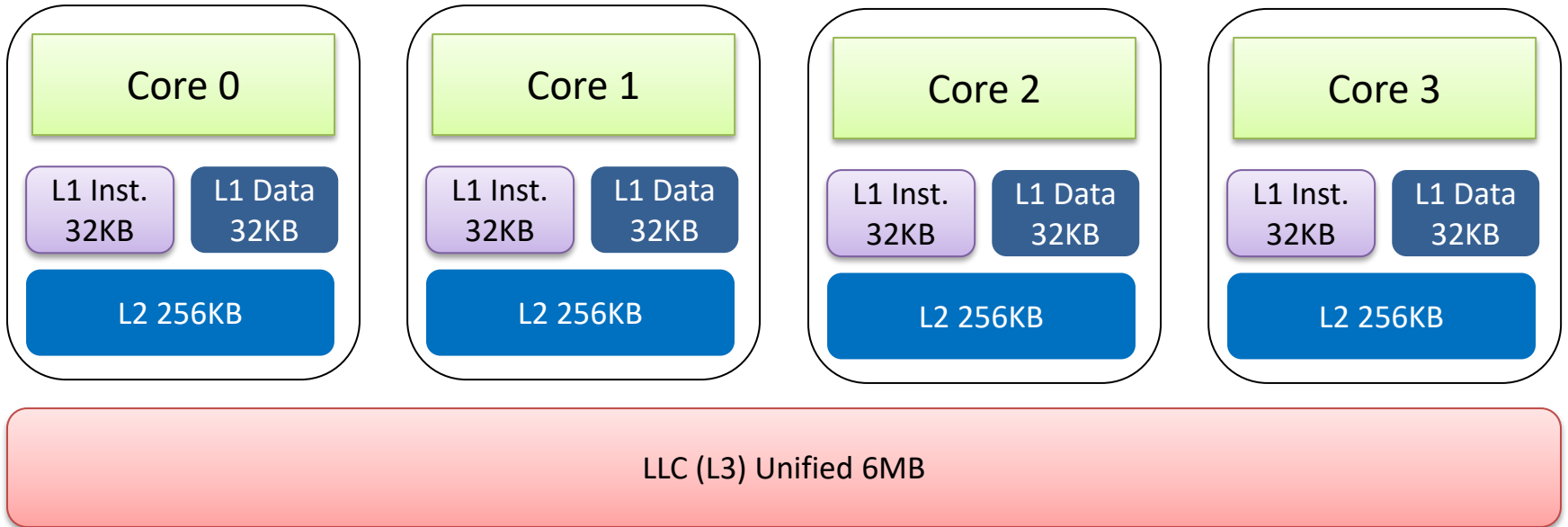
# Outline

- MRU is LRU
- LLC is not monolithic
- Use a crystal ball to help cache
- No, do not cache us
- Yes, you can control the cache

MRU is LRU

# Cache Inclusion

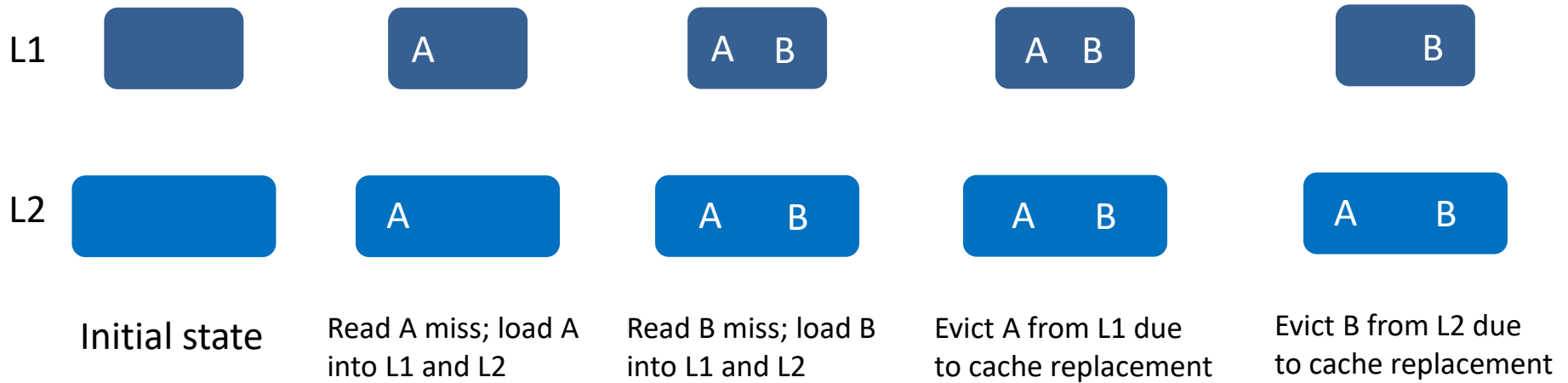
- Multilevel caches



Intel Ivy Bridge Cache Architecture (Core i5-3470)

If all blocks in the higher level cache are also present in the lower level cache, then the lower level cache is said to be **inclusive** of the higher level cache.

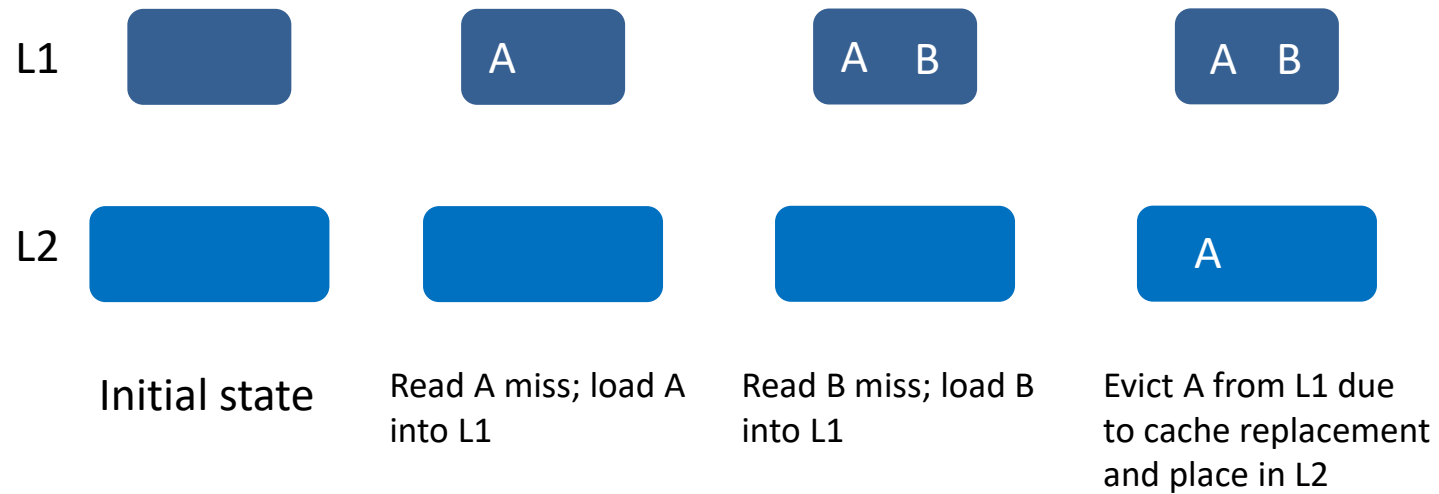
# Inclusive



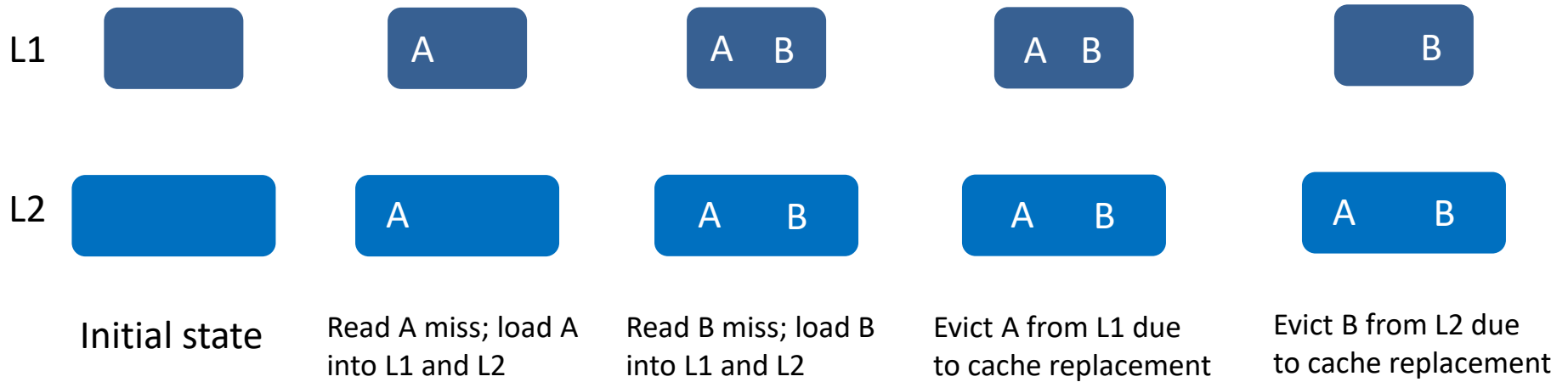
Back  
invalidation



# Exclusive



# Non-inclusive



# Real-world CPUs

- Intel Processors
  - Sandy bridge, inclusive
  - Haswell, inclusive
  - Skylake-S, inclusive
  - Skylake-X, *non-inclusive*
- ARM Processors
  - ARMv7, non-inclusive
  - ARMv8, non-inclusive

# Inclusive, or not?

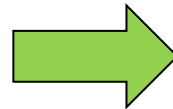
- Inclusive cache eases coherence
  - Updating a cache block in L1 entails an update in inclusive LLC.
  - A non-inclusive LLC, say L2 cache, which needs to evict a block, **must** ask L1 cache if it has the block, because such information is not present in LLC.
- Non-inclusive cache yields higher performance though, why?
  - No back invalidation
  - More data can be cached

# 'Sneaky' LRU for Inclusive Cache

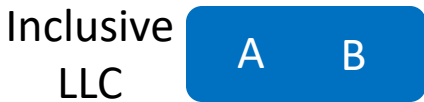
CPU  
Core



A is frequently used



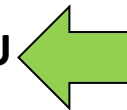
A is frequently hit in L1 cache. It is **MRU** in L1 cache.



A is evicted for replacement, in both L1 and L2



In LLC, A is **LRU**



In LLC, A is not frequently hit

As a result, MRU block that should be retained might be evicted, which causes performance penalty.

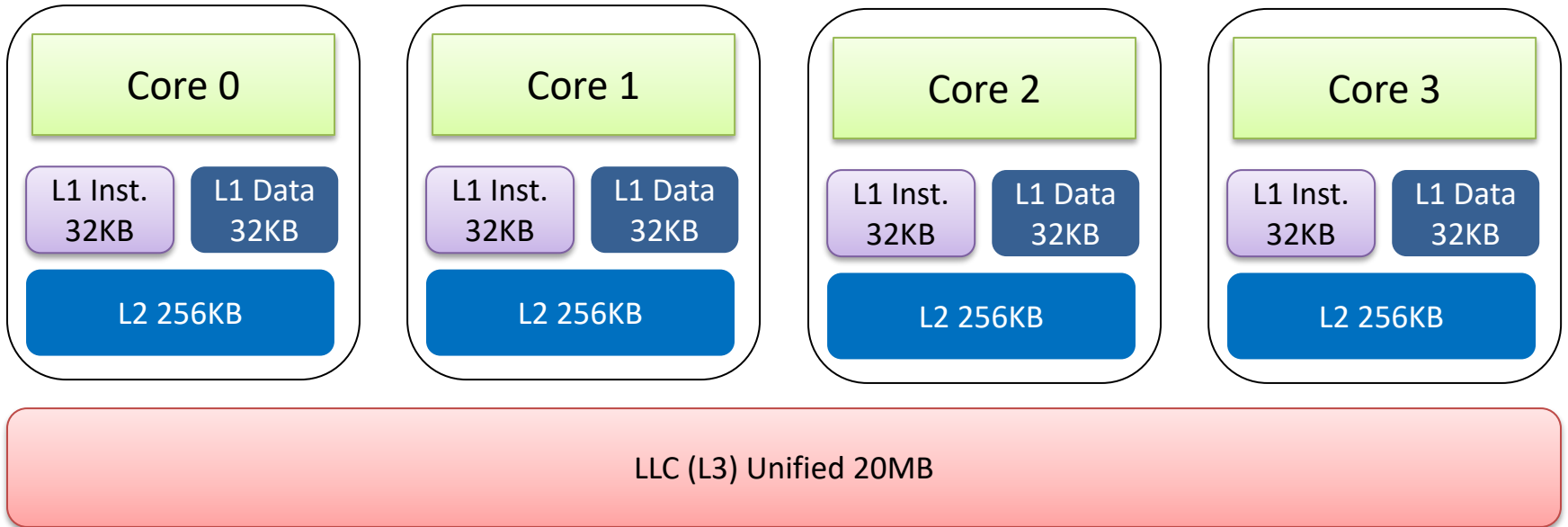
What if LLC is non-inclusive?

Should you be interested, you can click <https://doi.org/10.1109/MICRO.2010.52> to read the related research paper for details.

LLC is not monolithic

# LLC is not monolithic

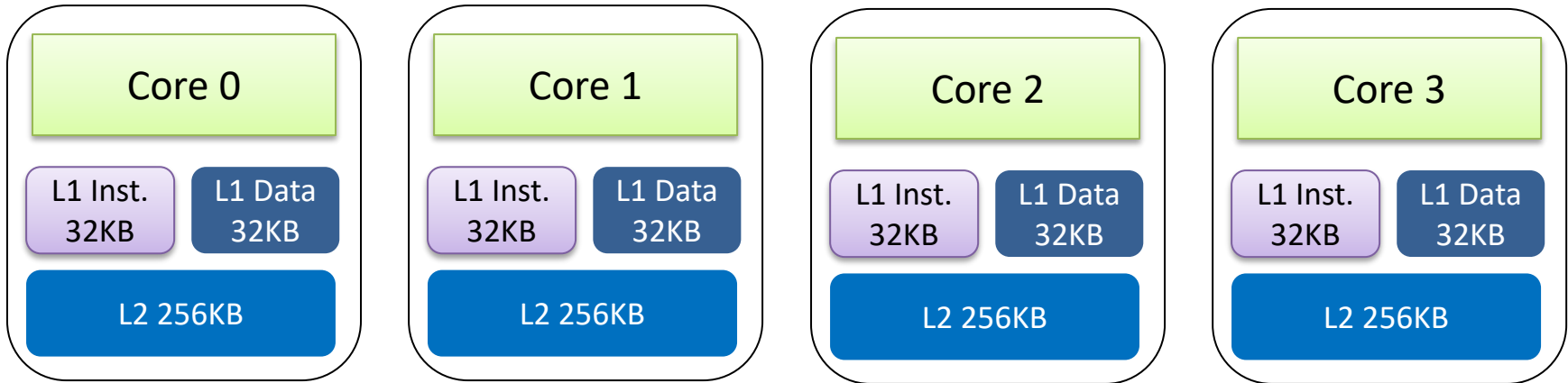
Intel® Xeon® Processor E5-2667 v3



Previously, it's considered that, to CPU cores, LLC is monolithic. No matter where a cache block in the LLC, a core would load it into private L2 and L1 cache with **the same** time cost.

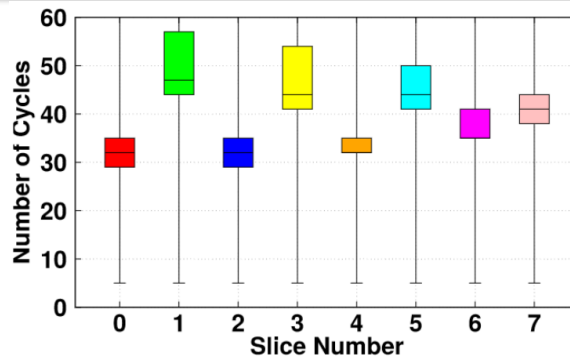
# LLC is fine-grained

Intel® Xeon® Processor E5-2667 v3

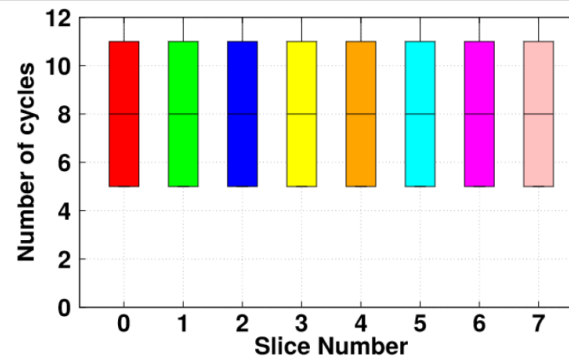


A

A



(a) Read.



(b) Write.



# Slice-aware memory management

- The idea seems simple
  - Put your data closer to your program (core)
- But it not *EASY* to do so
  - Cache management is undocumented, not to mention even fine-grained slices
  - Researchers did a lot of efforts
    - Click <https://doi.org/10.1145/3302424.3303977> for details
    - They managed to improve the average performance by 12.2% for GET operations of a key-value store.
    - 12.2% is a lot, if you consider the huge transactions every day for Taobao and JD

Use a crystal ball to help cache

Prefetch

# Outline of Prefetching

- Why prefetch? Why could/does it work?
- The four questions
  - What (to prefetch), when, where, how
- Software prefetching
- Hardware prefetching

# Prefetching

- Idea: Fetch the data before it is needed (i.e., *pre-fetch*) by the program
- Why?
  - Memory latency is high. If we can prefetch **accurately** and **early enough**, we can reduce/eliminate that latency.
  - Can eliminate **compulsory cache misses**
  - Can it eliminate all cache misses? Capacity, conflict?
- Involves predicting **which address** will be needed in the future
  - Works if programs have predictable miss address patterns

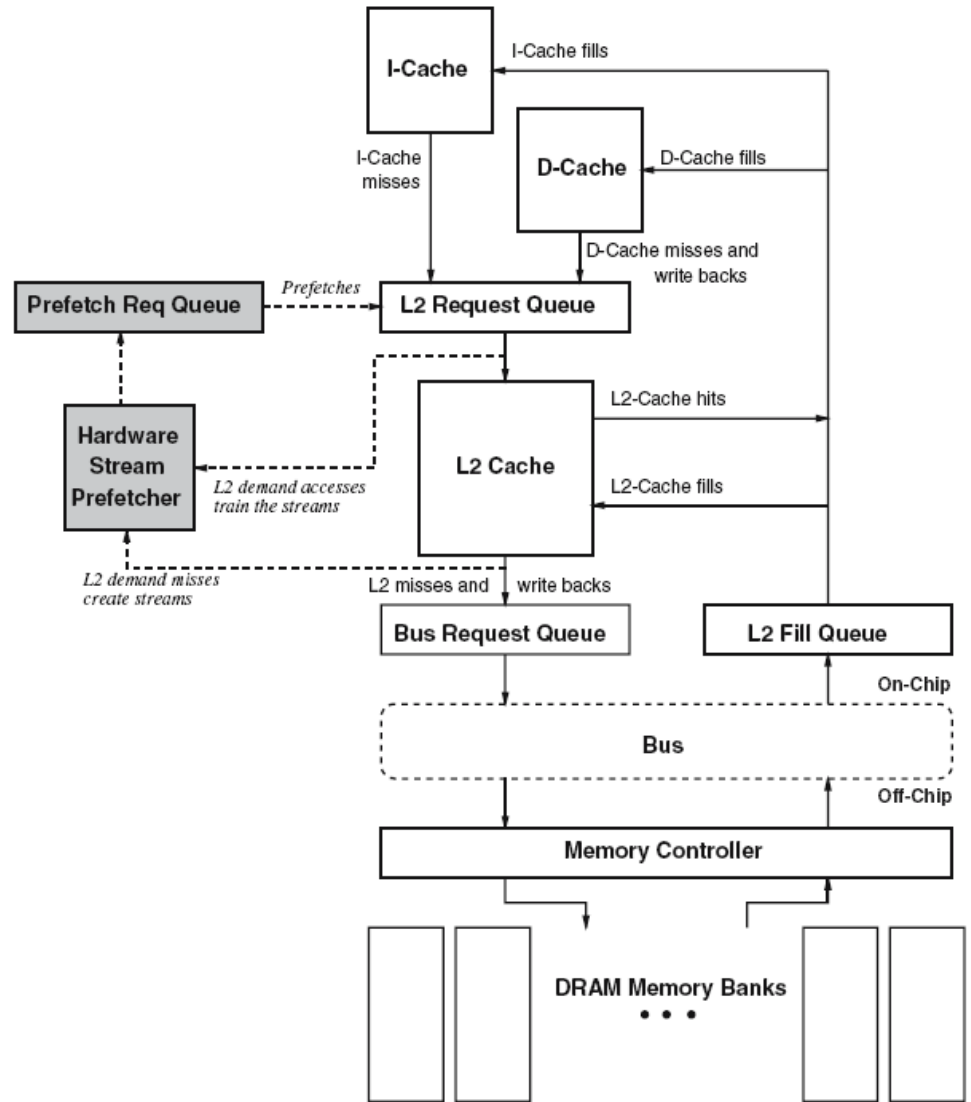
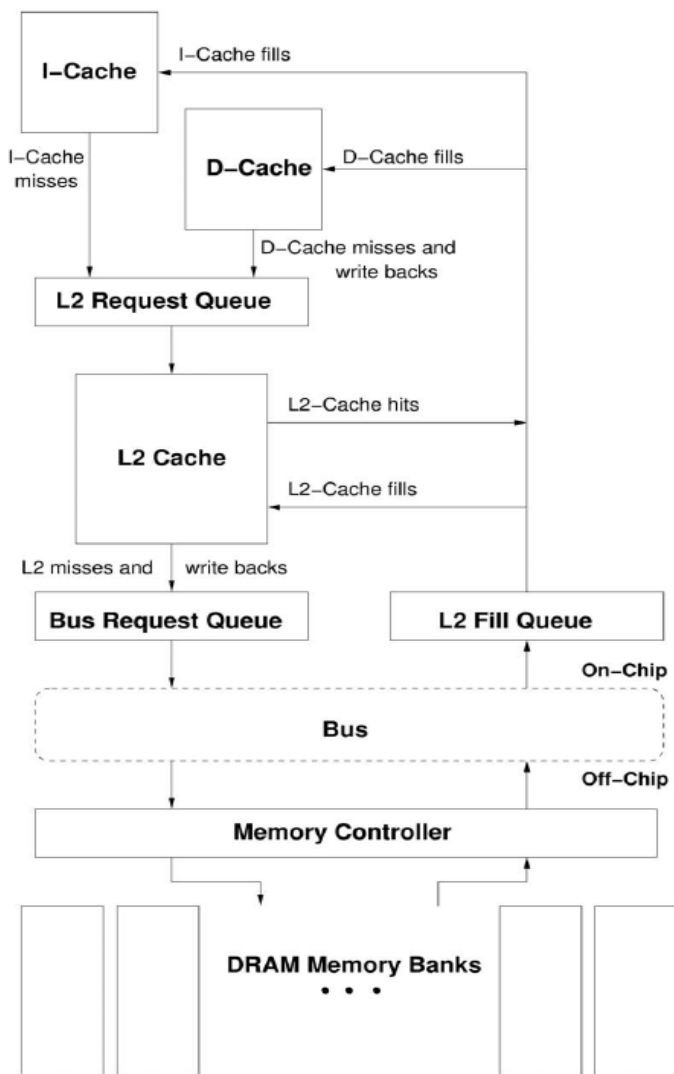
# Prefetching and Correctness

- Does a misprediction in prefetching affect correctness?
- No, prefetched data at a “mispredicted” address is simply not used
- There is no need for state recovery
  - In contrast to branch misprediction or value misprediction

# Basics

- In modern systems, prefetching is usually done in **cache block granularity**
- Prefetching is a technique that can reduce both
  - Miss rate
  - Miss latency
- Prefetching can be done by
  - hardware
  - compiler
  - programmer

# How a HW Prefetcher Fits in the Memory System



# Prefetching: The Four Questions

- What
  - **What** addresses to prefetch
- When
  - **When** to initiate a prefetch request
- Where
  - **Where** to place the prefetched data
- How
  - Software, hardware, execution-based, cooperative



# Challenges in Prefetching: What

- **What** addresses to prefetch
  - Prefetching useless data wastes resources
    - Memory bandwidth
    - Cache or prefetch buffer space
    - Energy consumption
    - These could all be utilized by demand requests or more accurate prefetch requests
  - **Accurate** prediction of addresses to prefetch is important
    - Prefetch accuracy = used prefetches / sent prefetches
- **How do we know what to prefetch**
  - Predict based on past access patterns
  - Use the compiler's knowledge of data structures
- **Prefetching algorithm** determines what to prefetch

# Challenges in Prefetching: When

- **When** to initiate a prefetch request
  - Prefetching too early
    - Prefetched data might not be used before it is evicted from storage
  - Prefetching too late
    - Might not hide the whole memory latency
- When a data item is prefetched affects the **timeliness** of the prefetcher
- Prefetcher can be made more timely by
  - Making it more **aggressive**: try to stay far ahead of the processor's access stream (hardware)
  - Moving the **prefetch instructions earlier in the code** (software)

# Challenges in Prefetching: Where (I)

- **Where** to place the prefetched data
  - In cache
    - + Simple design, no need for separate buffers
    - Can evict useful demand data → cache pollution
  - In a separate **prefetch buffer**
    - + Demand data protected from prefetches → no cache pollution
    - More complex memory system design
      - Where to place the prefetch buffer
      - When to access the prefetch buffer (parallel vs. serial with cache)
      - When to move the data from the prefetch buffer to cache
      - How to size the prefetch buffer
      - Keeping the prefetch buffer coherent
- Many modern systems place prefetched data into the cache
  - Intel Pentium 4, Core2' s, AMD systems, IBM POWER4,5,6, ...

# Challenges in Prefetching: Where (II)

- **Which level of cache** to prefetch into?
  - Memory to L2, memory to L1. **Advantages/disadvantages?**
  - L2 to L1? (**a separate prefetcher between levels**)
- **Where** to place the prefetched data in the cache?
  - Do we treat prefetched blocks the **same as demand-fetched blocks**?
  - Prefetched blocks are not known to be needed
    - With LRU, a demand block is placed into the MRU position
- Do we skew the replacement policy such that it favors the demand-fetched blocks?
  - E.g., place all prefetches into the LRU position in a way?

# Challenges in Prefetching: Where (III)

- **Where** to place the hardware prefetcher in the memory hierarchy?
  - In other words, what access patterns does the prefetcher see?
  - L1 hits and misses
  - L1 misses only
  - L2 misses only
- Seeing a more complete access pattern:
  - + Potentially better **accuracy** and **coverage** in prefetching
  - Prefetcher needs to examine more requests (bandwidth intensive, more ports into the prefetcher?)

# Challenges in Prefetching: How

- **Software** prefetching
  - ISA provides prefetch instructions
  - Programmer or compiler inserts prefetch instructions (effort)
  - Usually works well only for “regular access patterns”
- **Hardware** prefetching
  - Hardware monitors processor accesses
  - Memorizes or finds patterns/strides
  - Generates prefetch addresses automatically
- **Execution-based** prefetchers
  - A “thread” is executed to prefetch data for the main program
  - Can be generated by either software/programmer or hardware

# Software Prefetching (I)

- Idea: Compiler/programmer places prefetch instructions into appropriate places in code
- Mowry et al., “Design and Evaluation of a Compiler Algorithm for Prefetching,” ASPLOS 1992.
- Prefetch instructions prefetch data into caches
- Compiler or programmer can insert such instructions into the program

# x86 PREFETCH Instruction

## PREFETCH $h$ —Prefetch Data Into Caches

Opcode	Instruction	64-Bit Mode	Compat/ Leg Mode	Description
OF 18 /1	PREFETCHT0 $m8$	Valid	Valid	Move data from $m8$ closer to the processor using T0 hint.
OF 18 /2	PREFETCHT1 $m8$	Valid	Valid	Move data from $m8$ closer to the processor using T1 hint.
OF 18 /3	PREFETCHT2 $m8$	Valid	Valid	Move data from $m8$ closer to the processor using T2 hint.
OF 18 /0	PREFETCHNTA $m8$	Valid	Valid	Move data from $m8$ closer to the processor using NTA hint.

### Description


Fetches the line of data from memory that contains the byte specified with the source operand to a location in the cache hierarchy specified by a locality hint:

- T0 (temporal data)—prefetch data into all levels of the cache hierarchy.
  - Pentium III processor—1st- or 2nd-level cache.
  - Pentium 4 and Intel Xeon processors—2nd-level cache.
- T1 (temporal data with respect to first level cache)—prefetch data into level 2 cache and higher.
  - Pentium III processor—2nd-level cache.
  - Pentium 4 and Intel Xeon processors—2nd-level cache.
- T2 (temporal data with respect to second level cache)—prefetch data into level 2 cache and higher.
  - Pentium III processor—2nd-level cache.
  - Pentium 4 and Intel Xeon processors—2nd-level cache.
- NTA (non-temporal data with respect to all cache levels)—prefetch data into non-temporal cache structure and into a location close to the processor, minimizing cache pollution.
  - Pentium III processor—1st-level cache
  - Pentium 4 and Intel Xeon processors—2nd-level cache

microarchitecture  
dependent  
specification



different instructions  
for different cache  
levels



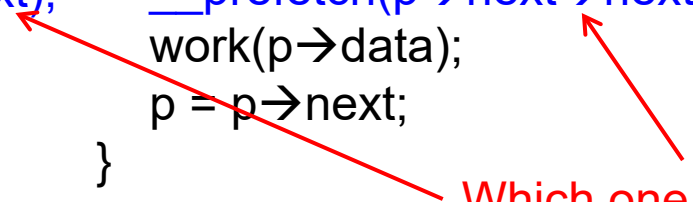


# Software Prefetching (II)

```
for (i=0; i<N; i++) {
  __prefetch(a[i+8]);
  __prefetch(b[i+8]);
  sum += a[i]*b[i];
}

while (p) {
  __prefetch(p->next);
  work(p->data);
  p = p->next;
}

while (p) {
  __prefetch(p->next->next->next);
  work(p->data);
  p = p->next;
}
```



Which one is better?

- Can work for very regular array-based access patterns. Issues:
  - Prefetch instructions take up processing/execution bandwidth
  - How early to prefetch? Determining this is difficult
    - Prefetch distance depends on hardware implementation (memory latency, cache size, time between loop iterations) → portability?
    - Going too far back in code reduces accuracy (branches in between)
  - Need “special” prefetch instructions in ISA?
    - Alpha load into register 31 treated as prefetch (r31==0)
    - PowerPC *dcbt* (data cache block touch) instruction
  - Not easy to do for pointer-based data structures

# Software Prefetching (III)

- Where should a compiler insert prefetches?
  - Prefetch for every load access?
    - Too bandwidth intensive (both memory and execution bandwidth)
  - Profile the code and determine loads that are likely to miss
    - What if profile input set is not representative?
  - How far ahead before the miss should the prefetch be inserted?
    - Profile and determine probability of use for various prefetch distances from the miss
      - What if profile input set is not representative?
      - Usually need to insert a prefetch far in advance to cover 100s of cycles of main memory latency → reduced accuracy

# Hardware Prefetching (I)

- Idea: Specialized hardware observes load/store access patterns and prefetches data based on past access behavior
- Tradeoffs:
  - + Can be tuned to system implementation
  - + Does not waste instruction execution bandwidth
  - More hardware complexity to detect patterns
    - Software can be more efficient in some cases

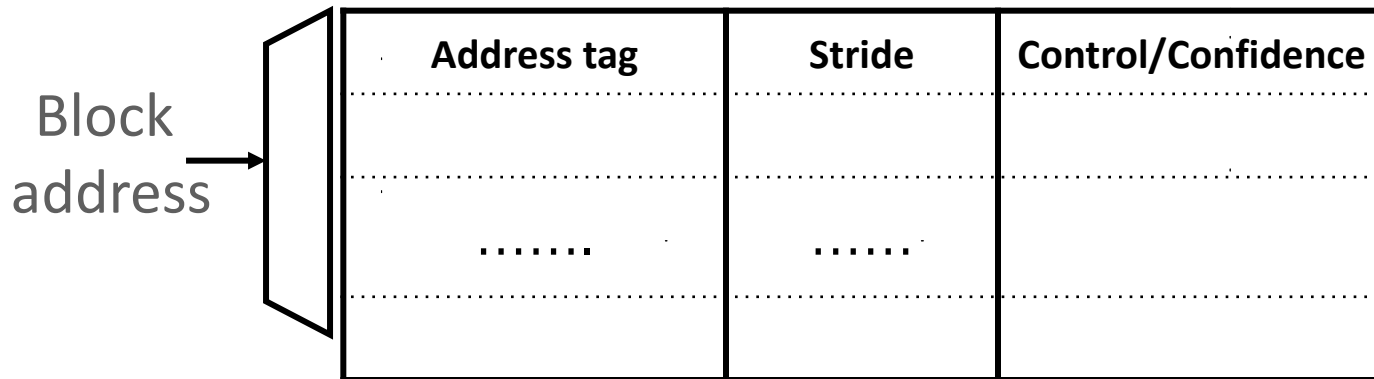
# Next-Line Prefetchers

- Simplest form of hardware prefetching: always prefetch next N cache lines after a demand access (or a demand miss)
  - Next-line prefetcher (or next sequential prefetcher)
  - Tradeoffs:
    - + Simple to implement. No need for sophisticated pattern detection
    - + Works well for sequential/streaming access patterns (instructions?)
    - Can waste bandwidth with irregular patterns
    - And, even regular patterns:
      - What if the program is traversing memory from higher to lower addresses?
      - Also prefetch “previous” N cache lines?

# Stride Prefetchers

- Two kinds
  - Instruction program counter (PC) based
  - Cache block address based
- Instruction based:
  - Baer and Chen, “**An effective on-chip preloading scheme to reduce data access penalty,**” SC 1991.
  - Idea:
    - Record the distance between the memory addresses referenced by a load instruction (i.e. stride of the load) as well as the last address referenced by the load
    - Next time the same load instruction is fetched, prefetch **last address + stride**

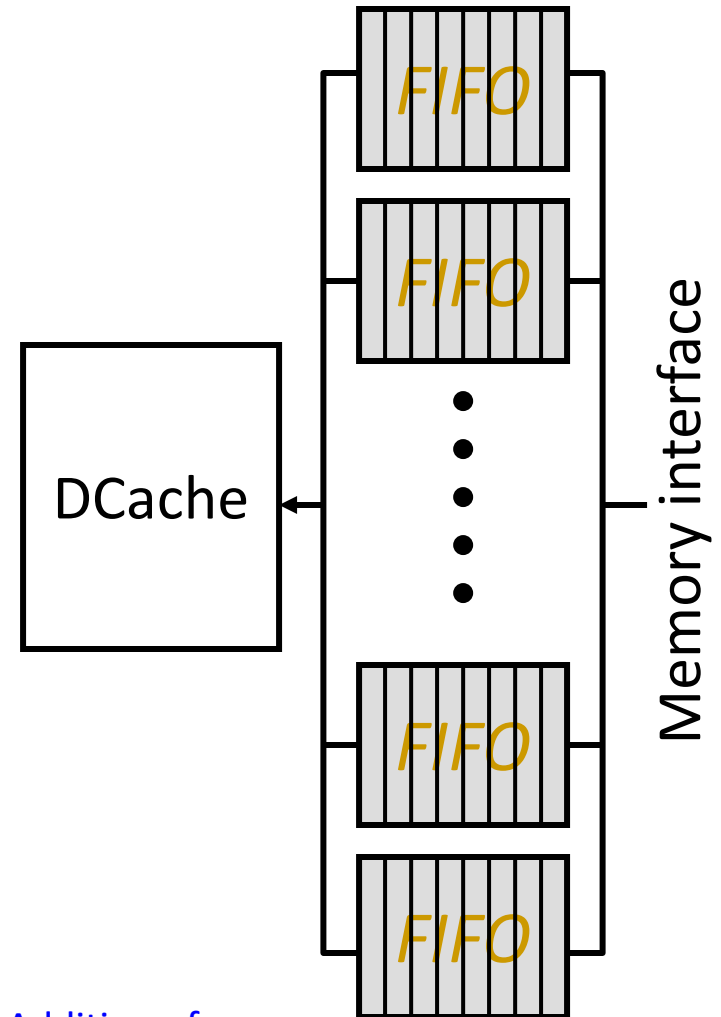
# Cache-Block Address Based Stride Prefetching



- Can detect
  - $A, A+N, A+2N, A+3N, \dots$
  - **Stream buffers** are a special case of cache block address based stride prefetching where  $N = 1$

# Stream Buffers (Jouppi, ISCA 1990)

- Each stream buffer holds one stream of sequentially prefetched cache lines
- On a load miss check the head of all stream buffers for an address match
  - if hit, pop the entry from FIFO, update the cache with data
  - if not, allocate a new stream buffer to the new miss address (may have to recycle a stream buffer following LRU policy)
- Stream buffer FIFOs are continuously topped-off with subsequent cache lines whenever there is room and the bus is not busy



Jouppi, "Improving Direct-Mapped Cache Performance by the Addition of a Small Fully-Associative Cache and Prefetch Buffers," ISCA 1990.

No, do not cache us



# Why do we use cache?

- The data to be used will be reused soon
- But there is some data for which caching is not that useful
  - So-called *streaming* data
  - e.g., larger matrices filled but used much later
- Caching such data pollutes the cache

# Non-temporal load & store

- Intel
  - MOVNTDQA: Load Double Quadword Non-Temporal Aligned Hint
  - MOVNTDQ: Store Double Quadword Using Non-Temporal Hint
- ARMv8
  - LDNP
  - STNP

Note that, such instructions only give a **hint** to the memory system that caching is not useful for this data

# When you can use them?

- The data is unlikely to be used soon
- The data is large
- Examples
  - Logs (journals)
  - In file systems and databases, logs are used for recovery/retrieval

Yes, you can control the cache

Scratchpad Memory

# Scratchpad Memory

- Strictly speaking, scratchpad memory (SPM) is **not** cache
  - Widely used in embedded systems
  - On-chip SRAM, like cache, close to ALU
  - Software controlled: software decides what data sections to be placed in SPM
    - By the programmer or the compiler before running
  - Memory-mapped to a predefined address range

# Why SPM?

- To control the execution time
  - More predictable than hardware-controlled cache
  - Especially for WCET (worst-case execution time)
- With reduced area and energy consumptions
  - More space- and energy-efficient

# Conclusion

- There are many interesting facts of CPU cache
- To make the best of cache can boost your program's performance!

# Quiz for prefetch accuracy

Piazza: "Video Lecture 23 Advanced Cache"

- What is the hardware prefetch accuracy if access stride = 1 and  $N = 2$ ?
  - A. 0%
  - B. 25%
  - C. 50%
  - D. 75%
  - E. 100%