

Technology Trends

آخر 40 سنة

- Transistors

- كدها فادرين زخيفت لحد حسب
moore's law

- تقريباً linear relation-ship

- بس الخرافة قاعد كذا بتاخر وقت لحد لكرير .

- علنا احنا زحمت خصال capacity بس ميس قادرين

Dark Silicon

نستعمل لكم نفس الوقت !

The amount of circuitry of an integrated circuit that cannot be powered on at the nominal operating voltage for a given Thermal Design Power (TDP) constraint

cluster التي على الجوالان دصا
مستعمل لارج حصة وباري يكون حياي ، وباري
مفنية ، او مستعمل on/off

- single-thread performance

ميس انتقلنا من ال single core الى multi core

- Frequency and Typical power

- كان توصل كذا 5GHz ، ولان بسط core ولا

بسبب لفرزك ← ال transistors المكونين power

Dynamic Power : بتقدر تفرج ال voltage وال frequency

فانزير ال F ، ال power تزيه

ال- Frequency توقف تست على أساس تقريبا لأكثر من مرتبة بال power والتبريد

كما حوا بالحوال voltage تست بال frequency: 2000 → 70/5

ال- P [تقريبا] تست بال P

quadratic تست بال

تست بال linear

ال- P

تست بال على ال- TTL ← العلاقة التي تست بال يكون فيها

logic 1 أو logic 0
+ve voltage -ve voltage

ال- ال- voltage التي تست بال على أساس ال- 20 تست بال - 1 volt

↳ range 0.9 → 1.1 (logic 1)

" 0 → 0.2 (logic 0)

بشكل قادرين نلاحظ هنا أن ال- تست بال على أساس ال- ranges

← وقتوا أيضا بال voltage تست بال، وكذا على frequency تست بال

قادرين نلاحظ هنا أن ال- تست بال على أساس ال- noise يولد هنا

مضطربة تست بال، تست بال، تست بال frequency

hierarchical.

Quick overview:-

احنا وقتنا بمشكلة ، احنا اصغرنا بال Thermal effects ، قولنا ان voltage قدرنا بقصر ، زدنا ان frequency قدرنا بقصر ، لكن مقدارنا اننا لغيرنا ، الحد الذي بيقتد على صريح ان voltage ، وال frequency ، وعلنا موجهة صحتنا قدرنا بقصر cooling

← دلستنا ان 2004 يطرح عننا ان multi-core processors

لأنه بيدينا performance ، ودينا progress وبتاكل لشغلنا خالص

- طرح ان intel dual processors / 2006

- " ان intel quad

- حبرنا 11 سنة شغالين عال i core ← i3 , i5 , i7 , i9

كفاي اننا صميم كانت موجودة قديماً ، لأن Moore's law فاساعدنا نتفكرها

لدهنا ، قدرنا بقصر ان cores على chip set ووجهة ، حلتنا بقصرها ،

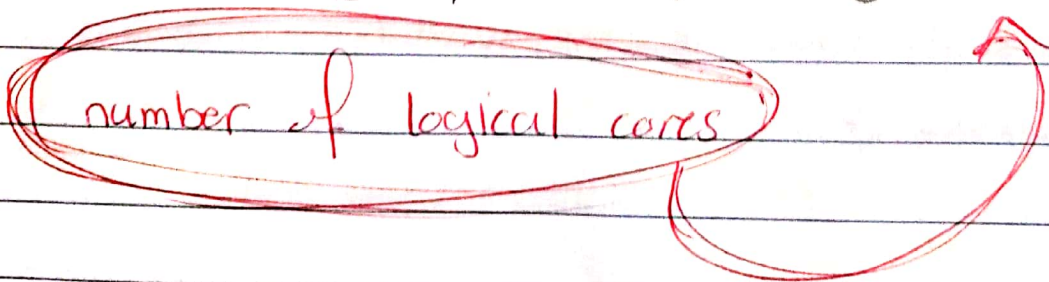
ان transistors

زمانه كانت تقدر بحركتها عننا ان core تانه يربطها مع chipset و CPU

كانوا تاولو زي chip ووجهة عليها ان core الرئيسي ، واولون حوايلها processors

ان تانه multi-core باعنا ، بتجي ، تانت arithmetic cores / GPU cores

ان built in حوايلها ان chipset ، وحصار بتكلم اننا سنالكه



number of logical cores

Notes:-

ال - power كاتبة تزيد ، حتى لانك كاتبة voltage و frequency
كناؤن voltage و frequency ، كاتبة بتزيد من تحقير عال
frequency ، حتى و على

النسبة ، عسرون
→

- عسرون نسبة ، كاتبة ال power بسبب مشاكل ، leakage current

كل ما زحيف ال transistors ، الالترونات عسرون → leakage current

ويتكهرب من ال path بسببها ، فلما زكرو ، بسبب ، ويتكرب من ال
أخذاً بعين الاعتبار انه سهل نكلم بال dynamic power بسبب نكلم بال

leakage current ، فوحطنا فرجة انه 50-60% من ال power ، (25%)

leakage

فبالنسبة حار عناه - Intel 3D Transistors / Mosfets
A field-effect transistor that has a thin layer of silicon oxide between the gate and the channel

بسببوا طبقات زكرو ، بسبب

وهلا ، مع ال leakage بسبب بسبب بسبب

Dennard Scaling:-

احدنا كنا نستخدم ال transistor ، فلما بدأنا نزيد طولها ،
كل ما اعل halving نستخدم ال transistor ، نبدأ بـ 1.4 صغرتنا ،
power density نقل للرج فرمونا

على مدى سنوات طويلة كانت هي ال rate

لقد قدرنا قدرتنا انزال ال voltage ،
مننا مع نبدأ من
مننا قدرتنا دفعنا ال voltage

لما اذا حيلنا بها ال trend ،
6000 درجة مئوية
processors

فبالتالي انجمننا ال multi-core ، ونشبت ال frequencies

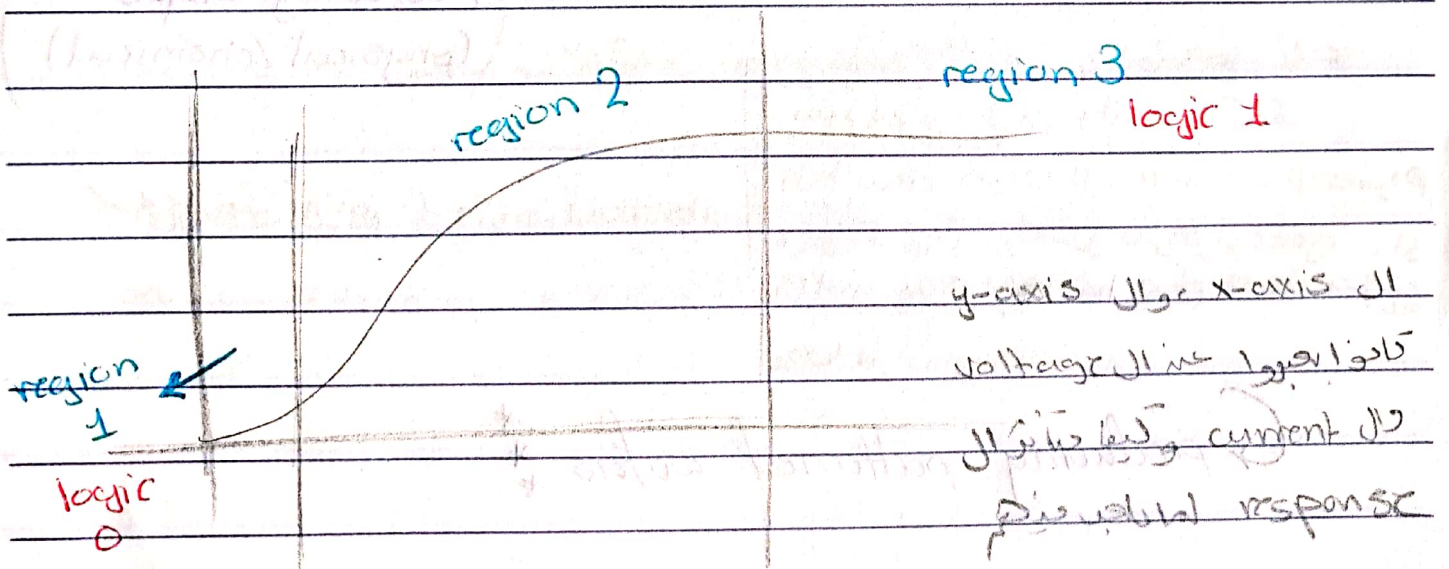
Slide (7)

- Basic mechanical tools →
- Mechanical computers → 1870 تقريباً
- computer systems → mechanical computing analogy
- كانت صيغتها
- وكانت تستخدم الحروف التي كانت قبل ال
- transistors → الى ال vacuum tubes

transistors gates ال

- VLSI → very large scale Integration → gates ال

* Response of transistor



computer systems ال region 2 بس 1 و 3

[Process of manufacturing CPUs] slide 18

- العملية تبدأ بالرجل ، صه ستواطي و بطاري مينة ، الرجل يكون فيه ستواطي و قلية

تصنيع و تنقية و حفر

silicon ingot ← تصنيع ← slices
slicing
slicer

clean rooms → الوان التي تصنع فيها wafer
Blank wafers ← طس حفر

(processing steps
(physical / chemical))

بفرد لبرازين ، ل ال circuits تبنة
ال circuits بنونها ل counters و حفر بفردهم
ل flip flops ، و حفر بفردهم ل gates ، و ال
gates بنفرد ل transistors بنفرد ل gates
بفرد ل transistors

← بنفرد ل لبرازين بال abstraction

→ producing patterned wafers ***

From patterned wafers :-

wafer testing

→ stage 1 of testing

يتم عمل اختبار على كل قطعة من شرائح الدوائر المتكاملة

يتم ذلك بوضع كل شريحة في جهاز اختبار ، وتطبيق إشارات الدخل على المنافذ

وتحليل إشارات الخرج من المنافذ ، وقياس سرعة العمل ووقت الاستجابة

ويتم ذلك باستخدام أجهزة اختبار متخصصة ، وتسمى هذه الأجهزة بـ wafer test equipment

التي تعمل على اختبار كل شريحة من شرائح الدوائر المتكاملة (IC) على سطحها

قبل إرسالها إلى العميل ، وذلك للتأكد من جودتها وخلوها من العيوب

والتأكد من أنها تعمل بشكل صحيح ، وذلك قبل إرسالها إلى العميل

processing

ال yield هو نسبة الشرائح التي تعمل بشكل صحيح على سطح الشريحة

↓

هذا يعني أن نسبة الشرائح التي تعمل بشكل صحيح على سطح الشريحة هي yield

→ Bond die to package (process for placing a chip on a package substrate)

Integrated circuit cost

كل ما زاد ال Die area ، كلما تزداد
الوقت ، كلفة تزداد

Measuring Execution Time

System time \equiv Elapsed time *

- احنا طابقت وقت تنفيذ كل لوقت الى بيلا ال OS على ال fast ال CPU time ، جرت ال overhead

لهيك قسمنا ال CPU time اجزتين :

System CPU time (2) User CPU time (1)

الكود الاصيل + الجزء

الوقت الى كينفذ البرنامج

اللي فاصرها ختيلها و

على كود البرنامج بيكون

المتعلق بال OS

اي بيكون آخر

لانه كل البرامج ختيلها طبقات ال OS

it is integrated

- Slide 1511

- Different programs are affected differently by CPU and system performance

two types of applications

- CPU bound \rightarrow

تخيل ان ال processing في ال CPU

- Memory bound \downarrow

ال shift ال ال disk

load/store ال ال برامج

RAM ال ال ال ال ال ال

disk ال

CPU clocking

- Gates (And, OR, Not) - they are only driven by their inputs.

- All our digital circuits are clock based

$$\frac{\text{clock period}}{\text{clock cycle}} = \frac{1}{\text{Frequency / rate}}$$

- clock is the main factor that affect the entire hardware

↓
كل ما زاد سرعة ال clock يزيد ال performance بافتراض
ثبات ال circuit design

→ in slide 17: CPU time is the total time not the user CPU time or system CPU time.

- there is always trading off between the factors that can be used to improve the performance (slide 17)

IC and CPI

How the compiler affects IC?

کے ساتھ ساتھ compiler کی بیرونی اور داخلی instruction mix, instructions کے ساتھ ساتھ assembly کے ساتھ ساتھ

- The hardware affects IC depending on program, ISA, compiler, and the design of hardware itself.

→ The actual CPI is like 10

T_c : cycle period.

§ 1.7 // Power Trends

Summary of power trend from previous lectures:

- Dynamic power \checkmark

- Static power \checkmark

- problems of leakage current \checkmark

- بدأت كثافة ال Frequency عند حافة صغرى ، او بدون تغير كبير انزلت
فترة بعد بين حبات بنفس الطول تقريباً

- كانت تستخدمه انه يزيدوا عدد ال cores في هناك ، ليوم هذا اليوم

- مقارنة ال power بتغير على ① خرج ال voltage و ② frequency

- حبرنا دخلنا ال voltage في 5V TTL ← 3.3V ← 1.1

الانام هاي رويست قدرين انزلنا اكثر بسبب ال noise بينظر تادرين

نميز ال 0s وال 1s في بوجون اذاتنا لاخت ال 1 volt

- قلنا دلنا تغير ال frequency ، كما انزلنا في حيز ، ودرنا ال

transistor في حيز ، وكان ال power density كتابة ← denard scaling

كتابة بغير ال 4 ارب ال transistor نفس المساحة مثلاً

بس نفس الوقت استهلاك ال 7 (4x4) بغير ال 1/4

تأخرنا ال

power

- The power wall

- can't reduce voltage

- still can (because of moore's law) ~~more~~ many many transistors

- These lead to increasing power density (because denard law has fallen)

but: still can't remove more heat caused by this power (no enough cooling capability)

→ slide // 25:

the trend is not straight.

حين دكسيتي، دكر في ساجو في الكسيتي

Multiprocessors

عند ال parallel programming
← مسألة كانت إنه تفاصيل ال Hardware تكون دغيتة أحياناً عند
الجروج، يمكنه يكون كود ال parallel ولكن ال code فاجروا كم هي
core عند ال user، فان ال actual performance ولا تحسين فيه عز معروفوا
إله وعز معروفون بلفارم.

← ال parallel programming صفة لليوم، كانه بيظروا
تركيوع حل، مسألة قد فان تركيز على ال performance نفسه

عندنا ما يكر ال load على ال core - load balancing
الواحد، وبقاين فربا حين ← توزع ال balance على ال cores

- ال CPUs عندنا في النظام، انه في عندنا حاجة اسما ال
throttling ← لانه، تنزل ال frequency للذي، عند تنزل
استهلاك الطاقة

Dynamic Frequency Scaling

لم لذل ال load balancing مهم، حتى انقادي ال
Hot spots، فبالأني ال hardware ما
دغيتا لسرعة، وبالسرعة أتا ما استنر

تركيوع لسحل على
واحد core

← في فرق بين separate programs (تطبيقات منفصلة) cores
 مثلًا microsoft excel ، core في word ، core في microsoft
 وهكذا ، كما يمكننا اننا كودنا في parallel ، اننا كودنا
 نقسم البرامج لواجهات يستعمل على كل ال cores في وقت واحد
 وهذا هو موضوع ال parallel processing

← باختلاف ال parallel programming كيف اننا
 communication + synchronization في كل ال cores
 مع نفس ال data

Benchmarking → مقارنه

Benchmark → قياس

قياس كمي انه كذا

القياس احسن منه كذا

← كثير website نسا profiling نسا ال CPUs و ال benchmarks

~~the~~ The standard benchmark in industry (The SPEC)



measuring CPU
time

- What are we benchmarking?
→ CPU?!
→ Entire System?! (response time)

Benchmark

تقيس النسبة
 - اذا لبي اهتم Benchmark، لازم اني load صند البرامج تقيس الاداء للكل
 - اذا جيب بوي اقيس اداء ال CPU جيب، فلانم اختار ال
 وقتها تستغل عال CPU، وتشي صند عال memory وهيك
 CPU benchmark Apps، التي بتقوي

- programs used to measure performance:

Supposedly typical of actual workload. →
 - لازم يكونه
 نموذجي لحي العمل
 النحلي

ال tools التي تختارها من
 جيب تختارها حسبك AP صند
 على مستوى شغل ال كبر صند جيب

machine for
 Geominy
 هيكلا لو بدي تقيس
 ديما تستخدم tools الحاسباته منها

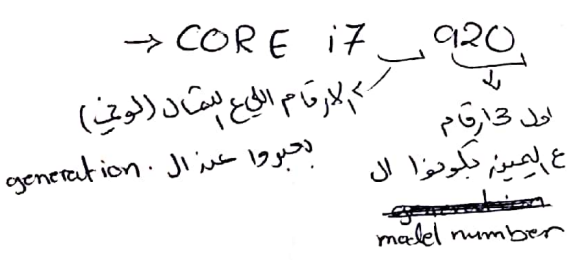
يحي هيكلا فالانم اختار tools
 ال user عمده حاج يستعملها

← ال apps المستعمه ك tools ال SPEC مقسومه دحين:
 ① مخصصه بال integers
 ② مخصصه بال Floating point.

- we have to rely on something, (a reference for comparison)

Case Study / Slide 130 → الجدول جيب عن البرامج التي مستعمه بالمعيار

CIUT: integer benchmark, CFP: Floating-point.



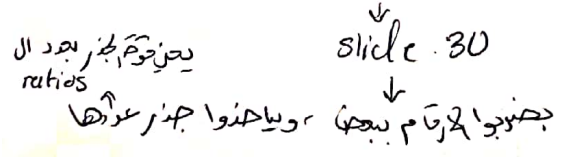
* هل هذا ال bunch of programs
 بيقر عن actual workload?
 له ناس مستعملين، فحاسب مستعملين

- SPEC ratio → جيب ال

$$\frac{\text{Reference time}}{\text{Execution time}}$$

to find a number that describes the measurement for the user:-

- use geometric mean not the arithmetic mean →
 لانه تباين
 بالقيم المتكويه



→ in the study case you'll take → product of ratios

TEC

عشان نفهم الحالة
والله ودهم لبيته.

- voltage
- capacitance
- frequency

static power
has leakage current. (2)
power distributed
for useful work or
unuseful work.

- on load 20% → استهلاك طاقة بدونها
استخر اتي

- static power موجودة دلياً
load. ← 5% حتى عال ان
نسبة كبيرة
total power consumption
هو تقال قيمها وقيمة طهر
لازم نسيم ان loads لاني

- كلما اطلع من load ل load اعلى ، كاتناسب بشكل حردى مع زيادة اوتيزال
Average power

[low power at IDLE] is a fallacy.

اجهزة google مثلاً دلياً مبرمجين ان 10% - 50% ، حسب الا عند سلكها يكون ان load ← 100%

static power نسبة ال
actual useful work اعلى من

- اجهزها ، انه نضم اجهزة فيها استهلاك طاقة تناسب حردياً مع ان load حتى نضف من هاي الحسبان

- desiring to design processors where power is directly proportional to load.

Pitfalls: Amdahl's Law $\rightarrow 1.10$

3 laws controlled computer design last-60 years:

- Moore's law is falling
- dinard scaling has fallen in 2004
- amdehl's law still affecting

→ "Make the common case fast"

$T_{affected} \rightarrow$ الزمن المتأثر

to know how much you improved.

$$\frac{T_{old} \text{ (long time)}}{T_{new} \text{ (short time)}} = \text{improvement (speedup)}$$

Pitfalls: MIPS as a metric

$$MIPS = \frac{\text{clock rate}}{CPI \times 10^6}$$

حفرجهوا بالحادثة (CPI) في حسابات HW
 وحاصلها، الحوارزمية، للوسائط كان تأثروا على ال CPI
 ليس عادلاً يتم قياسه بال CPU كانه

Example for Amalachi's

Program of 4 operations

Divide improving factor = 10 exec-time = 10

Multiply " " = 5 " " = 20

ADD/Multiply " " = 2 " " = 50

Other no improvement " " = 20

$$T_{\text{new}} = \frac{10}{10} + \frac{20}{5} + \frac{50}{2} + 20 = 1 + 4 + 25 + 20 = 50s$$

Chapter 4: The Processor

Pipeline: opposed to single cycle

hardware stages
 stages
 instruction
 HW
 instruction

registers
 pipes (pipes)
 stage
 propagation time for single cycle is big.

$$\text{Max Frequency} = \frac{1}{\text{time}}$$
 ← the min time to implement 1 inst

time
 frequency
 processor

Pipeline helped to solve these problems:

- parallelism
- pipes between each two stages, then the one stage becomes smaller, logic hardware of it is less, then propagation time of each stage is smaller

Higher Frequency
 path
 path

for this reason, ISA + internal design of computer, they affect frequency.

Five-Stage Pipeline

- 4 pipes

- why pipes?

- made from registers made from Flip-Flops

- Sequential circuits (Flip-Flops), they are memory returning circuits, we save data in them, then by these pipes (registers) we are returning information from the previous instruction as it goes to the next stage.

control signals وال operands وال opcode وال

to make this work

stages

Designs + Architectures

Harvard

for processors including embedded components or RISC-V

Von Neumann

Laptops + PCs

- one memory stick

per every thing

- managing data flow is the job of OS + internal system

- physically RAMs

in embedded systems

just

- memory sticks

for codes

- another memory

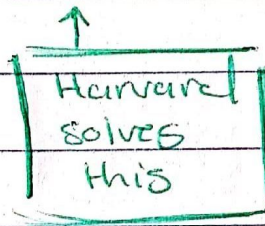
stick: data

that codes work

on

in von Neuman \rightarrow \leftarrow Harvard
von Neuman \rightarrow \leftarrow Harvard

\rightarrow Queuing problems appear for von Neuman, because one stick is doing the entire job, you have to wait.



- the big amount of wires in Harvard cause, electrical interference, causing a noise that affects the integrity of data
- in Harvard, each memory stick, can be made separately from different technology (SRAM, DRAM)
- Most von Neuman is DRAM most times.

lecture 06 min 14

Stage (3) : execute.

- the design in slide 5 is 64-bit
- 32 registers: length of each is 8 bytes : 64 bits

ALU:

- first operand is always coming from the instruction
- Opcode controls behaviour of ALU control.

Control signals

- RegWrite: decide to read from or write on register file stage (2)
- ALUSrc: ALU source; to choose one between imm, and inst. operand stage (3)

zero flag

- carry :

- half carry :

- overflow :

important flags
for designs. → most important

- Addsum: address sum.

unit ✓

calculating

target address

- when adding negative number
it is mostly for (for-loop)

- positive added number (forwarding)

Stage 4

- Branch control



zero flag
0 ↻

+

branch line signal
1 ↻



Branch instruction



you'll return

PCSrc with value = 1

going to

the mux of

PC to

take next or

previous instruction

← نبس هيل اناعرفنا اني بس اعلم Branching ، على ال stage

و فافترنا فنياً جالسنا ، و هيل كلمة يكونه كلمة instructions

على ال stages التي قبل الرابعة بالرئيس ، و بنسرحكم ← we flush every thing we took

into stages by making all control signals equal zero

just like we made empty instructions (bubble)

instructions (bubble)

stage (5) write Back
to write on register file

Question:-

assume in slide 5:-

stage 1: 10 ns

" 2: 15 ns

" 3: 50 ns

" 4: 15 ns

" 5: 5 ns

max frequency is $\frac{1}{50 ns}$

كل طول زمني لأطول stage فيه

بشكل الزمن يجب ان cycle



job of designers:

- to minimize the time

elapsed by the longest stage

مسألة صارت على انه لما نزل هذا

اصبحت pipes اكثر، وها ان

cost بتدوم على، لانه لما سيرة

sequential circuits من Flip-flops

بقي بدم clock وفضلنا جبرنا

interference على، وهذا كتر

- كذا كمنوتر، عدد ال stages التي قبلها

كتر بجزر ال branch صانحرو انا بكونه instructions

كثير اقل، وصرفت كتر با عالفاها، وودي flush

كل ال pipes

- another solution is to take the longest stage and divide it into many stages with less time, so that the longest elapsed time is less than the time elapsed by the original stage

Pipelined Control

- we need to determine the control signals, as soon as possible, directly after determining the instruction

→ using Opcode to determine control signals

Hazards:-

- structural hazards happens for the five-stage pipeline when the hardware tries to load many instructions at one time, then there will be competition on parts of RF and ALU, like when you load 4 instructions and you have 2 ALUs

→ a lot of instructions + few resources = delay and queuing

- Data hazards

like Read after write.

to try to read from a register you haven't write on it yet.

- Control hazards

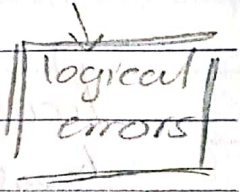
like late branch control.

The problem: ^{dependancy} ~~dependant~~ problem, sub inst / no problem, add inst needs to use the latest value of x2, x2 is found by sub inst, in ALU, and will be written on x2 reg in stage 5

Data Hazard in ALU instructions

Hardware level code analysis

and → decode stage
 sub inst
 pipeline



البرامج يكونها
 دورية نوقف
 ال ان
 كذا ما يكون
 ال x2
 صال sub

data تبع ال and وابقا ال ID

complete hold clock cycles
 الجزء لاول من ال pipeline ، وبدي ال sub
 ال باقي ال pipeline

Stall

and ال ID ، وال sub من ال memory
 ال رجوع ال RF ، كباي الكشاف يكون stall

contents of execute + memory stages
 → should set all control signals to 0
 nothing will respond to you (bubbles)

→ when add enters ID stage, sub will be in WB stage
 بقرا وبتاب x2 من نفس الكشافة على RF ، وها conflict
 فترعا ههنا كصير الكاشفة قبل القراءة

له عتانه كونر stalls وكل مشكلة بد ودهم ، بخلها بال forwarding

Forwarding Solution for Dependencies

- only for structural hazards (وحدیثاً و احیائاً جفتی توخت)

how to solve :- Before that, we have this solution

- clock has two edges for instructions in the case of add to solve the conflict

always write values on RF in first edge

" read " from " second "

to avoid the conflict



ال edge الاوّل بكتب

ال latest update

مع الثاني ان يكون

حيا ان القيمة

انكتب

This solves the conflict in between

sub, add.

slide 12: how many true dependencies?!

- and inst + or inst are true dependencies.

- add dependency is solved by hardware ~~forwarding~~ clocking

- sd has no dependencies.

total: 2 true dependencies.

Most likely in 5-stage pipeline,

we have the

2 subsequent

have the risk of

true data

dependency

Forwarding solution / slide 14

- instead of stall for and, or
- data is ready in EX stage of sub, why to wait until WB?
- do not read from RF for and, or, just take directly from where it is prepared
- other instructions find it in reg file

Forwarding unit | comparator ✓

① When to forward?

in next inst use an operand need to be evaluated in previous inst (1st operand or 2nd operand =

destination of previous inst)

inputs of FU:-

R1: 1st operand

R2: 2nd operand

EX/MEM. Register Rcd: destination of previous inst first previous one

MEM/WB. Register Rcd: " " " " second previous one.

outputs: slide 14

generates outputs it needs, forwarding them to the Mux (there must be control signals managing the mux: if forwarding is on, you take my forwarded data as an input for your ALU) otherwise you complete your path, with no need to any forwarded values, as there is no true dependency

solution of slide 15 in lecture 07 / min 22

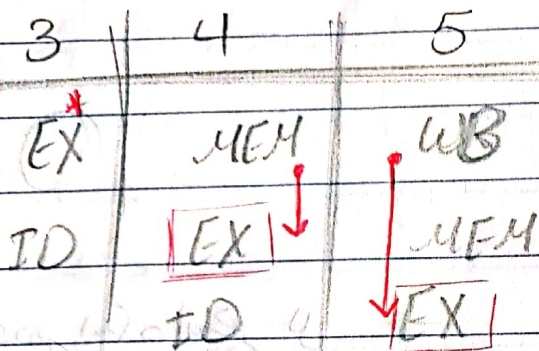
with forwarding design slide 15#

- X2 is ready in EX stage / cycle 3 / sub inst.
- ID of and won't read from RF the value of X2
- EX of and will read X2 by forwarding
- ID of or " " " " " "

قيمة X2 ، بتقرأها ال and وهي بال EX ، بالتالي ال sub تكون بال MEM ، ما بعد سلك مباشرة من ال EX بتعد Sub ، وسيتبع ال Forwarding Unit ، لا يقرأ ال X2 ال MEM بتعد sub same الطريقة هنا ، ال and وهي بال EX ، بتأخرتها ال sub بتعد ال or وهي بال ID ، ال sub بتعد ال or في هذا ال ID

← مدار فسر سلايد 14 لما ال FU بتقرأ
 ال RD مرة من ال EX/M
 وال M/WB بتكون ال sub

cycle 4 \downarrow تتم and
 cycle 5 \downarrow تتم or



→ we need X2 in EX, so we forward it to EX stage.

Forwarding works, and optimizing number of stalls if the instructions that are consecutive, are ~~arithmetic~~ arithmetic and logical, because their results, are prepared in EX stage. Unlike load-use operations and ~~data hazards~~.

slide 16

load: load from memory and write to RF

↓
 values are generated in MEM stage.

min 26 | left of

slide 17

assuming forwarding

- value of local will be prepared in MEM stage of ld

4 5 6

MEM

EX MEM

because of this, we are forced to

stay on ID stage

for sub inst for 2 cycles, until we can forward the loaded value from MEM of ld to EX of sub

لا نستطيع ان ننتقل الى next stage ال inst التي بعدها
لأنه يحتاج الى نفس ال cycle

cycle التي بعدها

عالم

حسب ذلك من stage

نحتاج الى بعدنا جاستنا ال WB ونحن نحتاجنا اعادتي بنقل
القيمة التي بعدها الى نفس ال cycle
عشاننا ال clocking

There is forwarding but with one bubble

Two ways to do local hazard detection.

①

ID/EX/MemRead
and
Rd \leftrightarrow RS1
Rd \leftrightarrow RS2

②

opcode of the inst
compared with
opcode of local
and

Rd \leftrightarrow RS1

Rd \leftrightarrow RS2

control signal

هو ال

slide 18

الوحدة التي تتخذ بها ال

load لا يقرأ ال

load سابق ال

reading from ال

memory.

طوكا (1) قبة ال

MemRead

Lcd ال

Hazard detection Unit

is found in D stage of

next instruction because

if hazard detected we

want to stall in D

of dependent inst.

Slide 14 // stall circuit → should have five NORs not 4, to be able to compare any possible register from 32 ones.

$$2^5 = 32$$

↓
comparing registers

truth table of NOR

RD and RS1, RS2

A B output

0 0 1

0 1 0

1 0 0

1 1 0

→ identical values = 1

otherwise = 0

A in this circuit → RS1 as example

B in this circuit → RS2 as example

←
checking the first operand.

There will be another exactly the same circuit to compare RS2 with RS1 (second operand)

→ The output of these two circuits will be both connected to an OR gate in order to say (yes there is dependency) if at least one of these circuits has an output 1

you can't ~~completely~~ stop combinational circuits \rightarrow (الذي يوافق ذلك) , but you can stop sequential circuits.

\downarrow
Flip-Flops \rightarrow let them keep previous values but not to update.

\downarrow
How?

First solution: stop the clock and set it to 0

\downarrow

مدخل ال clock نفسه ، كيفية ال clock
التي ال clock ، فلو دخلنا 1 بديلاً
input is and gate .
ال clock نفسه ، ولو دخلنا 0 تبطل
ال clock موجوده .

\downarrow

من ال clock

\rightarrow The second way !!

Way 2 do not update!

→ PC counter should keep the previous value, and never update.

→ every thing stored in the pipe between IF and ID will remain the same

لأنه لا يتم تحديث الـ clk على القيمة السابقة فالتوقف
وال flip-flop يكون disabled ، فكل
البيانات التي في الـ pipe
ستبقى كما هي stall

↑
من الـ update الـ register
على الـ value الـ الـ
register
بإزالة ← disable

Design of Register

D-Flip-Flop JK-Flip-flops

→ inputs (Q, \bar{Q}) or D

→ outputs

→ clk

→ set, reset

→ there will be enable / disable wire.

Data path with Hazard detection

Hazard detection unit

inputs: RS1, RS2, RD

+ ID/EX MemRead control signal to detect local hazard



إذا شرط الشرط الذي يتولد فيه خطأ
بإيقاف الـ control signals

التي هي

* First output: to the mux in ID stage that controls WB, Mem, EX control signals

→ control lines pass as usual

→ if hazard is detected, the selection lines forces the mux to pass the zero value and assign it to all control signals mentioned above

Two inputs of mux:

- 0

- control.

second output of Hazard detection unit:
- PCWrite; to disable if needed.

Third output,
- IF/ID Write; to disable pipe if
needed.

Solution of slide 25 in lec 07
min 44

example slide 26 in lec 07 min 49

In General what is happening now

making:

- registers of size 64-bit
- addresses of size 64-bit
- any thing loaded from the memory is of size 64-bit
- instructions depend on the architecture: 4 bytes or 8 bytes.

In slide 11:-

when add is in EX stage (will read from RF) and sub is in WB stage (will write on RF), and both are happening in the same time in the same cycle (This is a conflict)

This solved by clocking

- The job of rearranging codes is the job of an intelligent compiler, ~~the~~ it does filling of such an instruction that has no true dependency between two instructions, ~~that are depending on~~ where one of them depends on the other

↓ To solve load-use hazards, because other data hazards can be solved by forwarding.

In slide 7

→ the operations and calculation in order to determine we have a branch is actually taken is done in EX stage, using zero flag an output + calculating the next branch pc

- But the real action is going to be done in MEM stage, (we don't take action on branch is taken) except in MEM stage.

لأنه كما نرى هنا يكون داخل 3 instructions بطول ال pipe line وهذا هو ال Fall through ال sequential instructions
المرتبين ورا يوجد بعد ال branch [وهدا وهدا وقت] وقد يهرون
السابع

let assume that we have:-

$beq\ x1, x2, 16$
→ offset.

using design of slide 7:-

- instruction is fetched in IF stage
- registers are read in ID stage
- $(x1 - x2)$ is done in the ALU of EX stage producing a zero flag, for which the zero value means, $(x1 = x2)$
- then the branch is taken + next PC is found here by (AddSum circuit)
- producing zero flag as an input to the branch detector logic (And Gate), with control signal of branch in MEM stage **which is too late ***

Observation #1 update

$x1, x2$ are read in stage ID, so we can work on them as soon as ID stage.

Observation #2

PC address which is used to find target PC is detected in IF stage and is passed through all stages until it reaches AddSum circuit in EX stage.

Observation #3

to know if the inst is branch, you can do that in ID stage by (branch) signal, you are not forced to wait until MEM stage comes and use this signal.

Solution #1

you can rearrange circuits in a way you finish the job as soon as ID stage and protect 2 cycles from stalling.

Target address in branch inst. is calculated

by adding (current address + shifted literal by 2)

in slide 30

$$40 + (16 \times 2) = 40 + 32 = 72$$

↓
you first sign-extend this literal and then you multiply by 2, by doing shifting in hardware

This solution is represented in slide 32

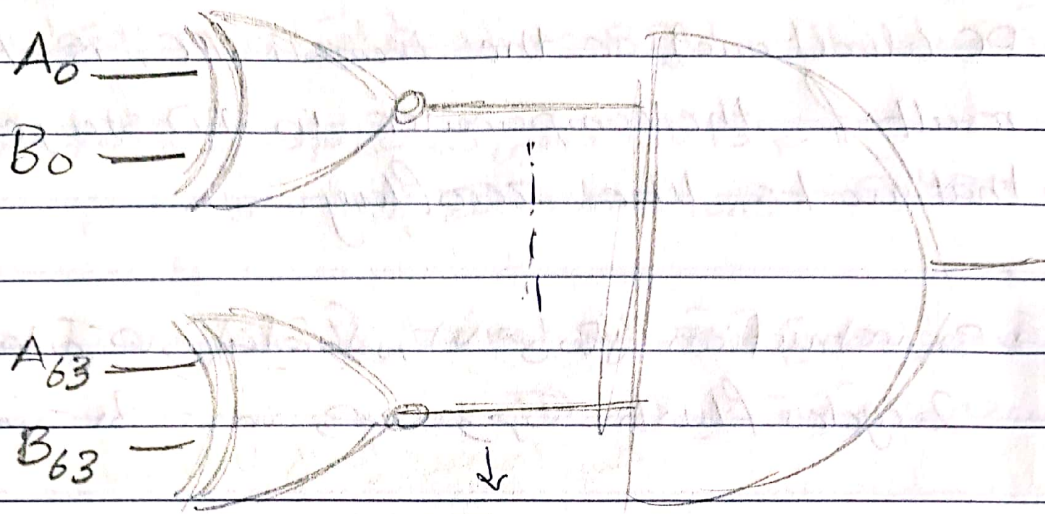
① - comparative circuits directly after RF to compare between $RS1, RS2$

The design of this comparative circuit

it is made from ~~XNOR~~ XNOR Gates

if $RS1 = A$, $RS2 = B$

Then:-



- you have 64 NORs here
- This comparator is like the stall circuit in slide 19.

But NOTE: AND gate can't have more than 8 inputs, then NORs are distributed to 8 ANDs then these ANDs are connected to one last AND

- This step detects if the branch is taken or not.

* (2) AddSum circuit that calculates the target PC, is put in ID stage not the EX

→ The shift circuit that multiplies the literal by 2 to find the offset is put with AddSum in ID after extending it in ImmGen.

③ The selection line of the mux in stage IF that decides what to take, the next PC (direct one), or the branch PC, is the result of the comparator in ID stage that works like zero flag.

بما اننا نستخدم instruction في 2 cycles flush

This solution is not optimal!

because the designs of these days are deep pipeline designs, and this solution can't help in them.

BPU: Branch Prediction Unit.

Solution #2

solution of RISC-V pipeline

طابقنا تصديرات ال BPU ربتنا وحدة اقتراح واحد فقط ونسبنا على

على اجراءه اذا كانه كحسي في على معرفة على بكمية ال instructions بكمية لإعداد افضل

على ان يكون انه يتم توقع ال branch بال IF
وذلك حسب قيمة حقيقة ال ID

giving not taken prediction

هون اصابعه حارة ال BPU دابة

36 سلاية

على اقتراح انه ال BPU دابة بظي ← not taken

فصل اجراءه قاعيه بنجل ال full through sequential instruction

لربنا تقري انه ال comparator ال ID ال على ال
انه ال RS1, RS2 كساحوا بظي ال branching

ال حال ال comparator ال ID ال على انه Taken
خلاقاً ال BPU هون اصابعه بنجل ال Flush ال inst

| | |
|---|-------|
| هذا الحل مسدداً بزيح حسب الوقت قدس ال branches واحد | cycle |
|---|-------|

This is called Predict Not Taken

The other design is Predict taken

improvement #1 on solution 2#

static branch prediction

Backward

Taken

لا تفرغوا

loops

cost

Forward

Not Taken

if = ال

else = ال

نا وخطي

wrong ال prediction

هوية تكون اعلى

not taken ال

taken ال

عدد لغات كتار

improvement #2 on solution 2#

Dynamic branch prediction

what is the history of this instruction behaviour?

- buffering:

عينة كثرين

Big BHT: cost / time problems
small BHT: conflicts problems.

Table size = $n * 2^k$ bits

n : n -bit counter

k : bits you take from the addresses

→ A finite state machine

- The one bit predictor does the job of updating data in BHT

مسئله این جدول

(عای آجز update له)

طرح مخالف بروز تعد

ال table سرعت و

اذا توقعات اختلاف

و سرعت (بفرض هلك

تكون و تبا تر كماله بال

comparator مع ال ID

های مسئله تبیین لها بکون →

for loop داخر

for loop



slide 41

assuming

→ inner loop will be executed

1000 times

→ outer loop will do so



million times you'll reach bee in the table.

2-bit Predictor

| | | | |
|------|---------|-------|---------|
| Deep | Shallow | Deep | Shallow |
| NT | NT | taken | taken |

نقل ما تم نقله مرة ثانية من NT إلى T
 shallow taken → Deep taken
 ← إذا خرجنا من NT، (الرجوع) إلى pattern أو قسماً
 نقل ال prediction

~~Improvement~~

Improvement on the whole design:-

- predict in IF and calculate target PC also how?

① First time you access bcc inst as example, it will be calculated in ID by AddSum, and then save this target address to the corresponding inst address in BHT in the same row with its 2-bit or 7-bit predictor value, to call this table BTB after adding the target PC
 ← هذا هو ال predictor في ID
 ← على ال ID نكتب ال target PC في ال BHT
 ← ال BHT في ال IF ال predicted PC

BHT: without target PC

BTB: with target PC

Some hardware divide the to 2 table

T1: inst address + bit-predictor

T2: " " + target PC

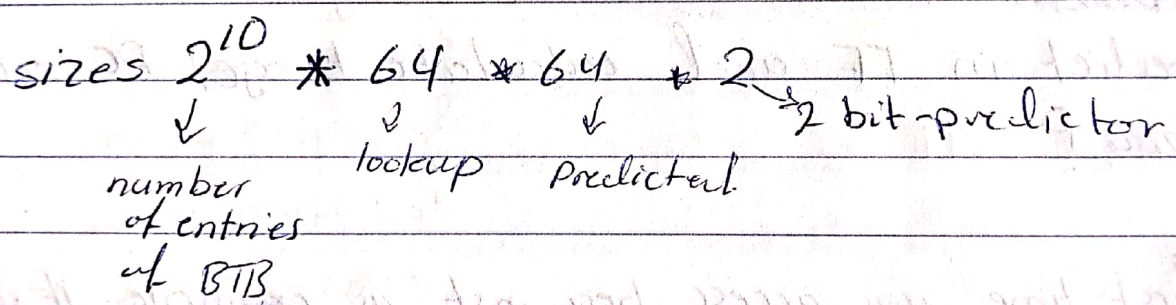
Another solution:
 is to ~~with~~ build an ALU or adder that
 calculates the target PC in IF instead
 of depending on the one in FD

IF, ID, W, E, M, S, P, D

Example:-

what is the size of BTB where

- $k=10$ (number of bits taken from addresses)
- R64 VI Architecture.
- Instructions of 64 bits.



length of table = 2^k only.

Exceptions and Interrupts : Control Flow Hazards

control Hazard on flow عطل في تدفق البرنامج
flush

التي يتم اعادة توجيهها الى البرنامج الصحيح

But **Exceptions are** → انسي مشكل في كبر داخل ال CPU نفسه
انما تنفيذ البرنامج

لهذا اعلى الاوقات يكون مشكلة واصحابنا ال exception يكون انه علامة
بال OS

what kind of problems?

① problem in IF stage :

ال system يكون في noise في مكان لا علاقة له بال inst
memory في ال hardware ال noise في flip bit في خرج ال address
بغير انسي مختلف

[WPU] noise protection unit

يعني inst bit مشكل ال noise في
السبب الوجود

اذا حدثت مشكلة في ال hardware
بوجود انذار

An exception in
fetch stage

if stops your code:
"you are trying to
access invalid
memory location
.... etc."

② problem in ID stage

undefined opcode

عندما تكون 7 bits opcode
التي ليست مخصصة
128 instructions

80 arch في 80 inst ← لكن
undefined

ال assembler or compiler
ال opcode غير

لكن لغرض الأسباب ال noise
عند هذا ال opcode
ال undefined

The 80 is an example

So in ID, this opcode is entering the control unit, that is designed with truth tables of logic to recognize only our 80 defined combinations, but now it's seeing something undefined and it's not designed to recognize it.

This rises an exception (ID exception)

③ problem in EX stage

→ two operands } → not 2's complement representation
→ unsigned

لا تفرق بين الـ 2's complement
والـ unsigned

range of it is

we designed overflow circuit. \leftarrow

(overflow)

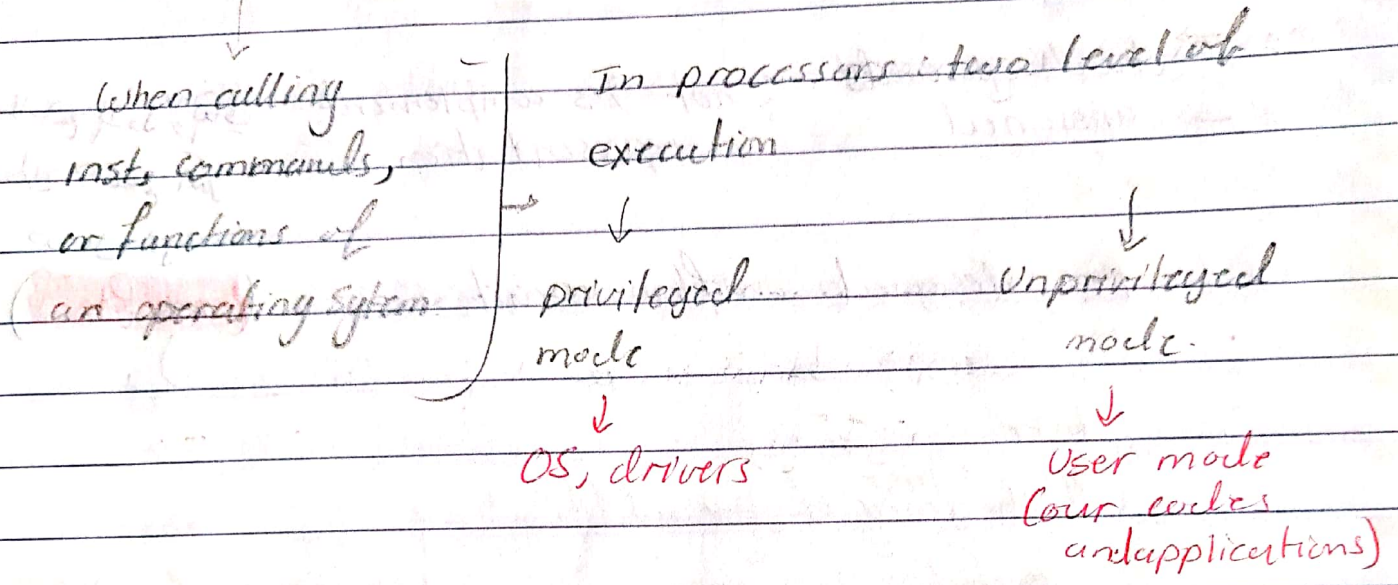
take:

- ① the last carry of the addition operation
 - ② the least sign bit of the result
- inputs for XOR
1: overflow
0: no overflow

overflow: exception

These exceptions are disabled in many arch, we enable them manually. (Embedded controllers and processors)

④ syscall | exception



which is OS

↓

the instruction that tries to communicate with HW component, will give a syscall → producing an exception that gives the handling of execution to OS to do the job and return results to you

we should have an intermediary component.

they can't communicate with HW components directly

And **Interrupts** are: → dealing with input-output

Flush circuit → Re-fetch from new place

الاضطراب interrupt سؤالي

عندما circuit في حالة اول سؤالي input في اول mouse سؤالي
في اول circuit سؤالي امسج ال processor ال سؤالي interrupt سؤالي اول
inst ال سؤالي اول ، و سؤالي ال instructions سؤالي كود ال mouse سؤالي

سؤالي

Flush circuit ↔ Re-fetch from new place ←

Way (1) # scenario

Question: we want to get the mouse instruction, after we received the interrupt, can we go directly?

→ interrupt is received,

→ we saved where to go back again (interrupted PC) in SEPC

→ saved the indication in SCAUSE

→ intending to go to the code of mouse, or to any code that can solve the problem if exists

how?

we don't know where are the addresses of instructions in memory (HW does not), but the OS does so!!

So: ISR: interrupt service routine or

IHS: " handler subroutine or

IHF: " " function

address of ← This function where we are directed to, after finishing steps of SEPC, SCAUSE

→ We go to handler, automatically, this function is simply an if-statement code

OS solution instructions ← SCAUSE

SEPC + طیب جس کی مسئلہ بہتر ہے کہ PC کی چیز تیار ہے

↓
problem کی وجہ سے

OS کی crash ہوگی
will throw an exception

could return?
restartable
Handler
could not return?
termination
happens.

Not All designs are alike

Way 2 # scenario

if-statement code اور address کی وجہ سے
exceptions & interrupts کی وجہ سے
address اور

no need to handler

faster
solution

vectorial
interrupts

Hard coding

Job of
programmers

to handle exceptions.

more complex
hardware and
higher cost

Not easy!

Exceptions in pipeline

على كذا جزء ال PC أو PC+4 حسب نوع instruction ال SEPC
↓
the next to the offending

مثلاً ① يصبح invalid opcode، فاستجابة ال instruction نفسه، ديو ياتي ال handler ليحل المشكلة فيه كذا بيان، جودي استجابة ال OS وهو يبرمج ال نهاية حالة هو ديو ال PC

② عند ال interrupts، احنا دينا قبله ال next instruction
لما نرجع، فنغير ال PC+4

[Improvement on the Pipeline]

slide 55

① تغيير ال signals

- signal IF flush on the mux before PC
- " ID " " " " in ID
- " EX " " " muxes in EX
- ال control unit ال flush ال و ال to decide per-WB signals
- MEM signals.

control signals هو ال control unit، ال overflow

لما درس ال design ديو ال (10) دقيقة 22 اسطوان ال mux
كبار ال اذ كانت ولقد.

الأجهزة المتصلة:

← لو حصلنا مشاكل بال stages التي قبل، راح عترة wires ال SEPC
SCAUSE قيمهم كبر ما توصل ال registers بتخدمهم بال ال EX stage

الأجهزة المتصلة:

عنوان ال IHS بيختره input عبر ال mux الي قبل ال PC

slide 55

Restartable exceptions

NonRestartable

Go back to complete

termination

program

← بالذات ايات احقة ال SEPC SCAUSE يكونوا موجودين بال كل stage
حتى ما رجع cycles لتقدر توصلهم زي المرات

Exceptions Example slide 57

Handler IHS

↓

one interrupt - vectored

→

IHS code instructions

- sel .. why store?

↓

Handler code is if-
statements

↓ context switching

تغيير البرنامج

to save on stack

دهفظ القيم بالذاكرة ال RF

بال ال MEM عنوانه لما نرجع منه على ما وقف، عايننا قيم ال inst لقيمة التي وقف
عندها حاجة

→ Handler

Store instructions

IF-statements

Local instructions → to restore :-

SEPC ← لأنه مجرد ترقية البيانات وان SEPC
تستغل الكود القديم First: ~~SEPC~~ old values
then: ~~SEPC~~ SEPC

← وطاقام الختار فحنا لمشكلة، بينما نوضح ان Flipped-bit ان
SEPC ← SEPC
منه 1 لـ 0، عشانه لو بيدي exception جبره ويوال OS بيقدر
ان SEPC هيك بيتوقف في SEPC، ما بيتوقف لمشكلة قديمة
ديتريبط فيها ويزرع علكانه غلط.

← بالذمة الي يعرف فيها مشكلتها - سجل clear flags، ان SEPC

* we'll lose three cycles in this example

Multiple instructions

Question: Could we have more than one exception?

yes, it could happen, but if so, multiple exceptions won't be in the same stage.

[slide 61] lecture 10 Min 32

- add instruction have exception in EX stage / cycle 3

- I3 (bad) : invalid opcode. دور على ال opcode
تتم بمرحلة bad
will be detected in ID

↓
* 2 exceptions in the same time / same cycle

↓
to solve them, they are ordered depending on priority in modern designs

- exceptions have higher priority than interrupts

↓
لا كما يتصور بان CPU

دائماً الاطوية لدى سي

دوران CPU

← لسانه اين built-in exceptions اول، اما ال interrupts اصنافهم

اولوياته بان ← software ← programmable ← او بتصميم اولوياته
عند مرتبة ال اسلاك
وتشكيهم

← في ريزاين ال RISC-V ، بتسجل ال priorities ال software

→ Solution is **precise**

they happened in the same time

but you solve the older in the sequence of code

(adder have higher priority here)

the operation

→ Fill SCAUSE with add bit

→ Fill SEPC with PC or PC+4

→ IHS solves add exception

→ return to I3

→ I3 will get stuck again in ID

→ exception detected.

→ Fill SCAUSE + SEPC

→ jump to IHS to solve exception

High Cost

Solution

+

Overhead

→ Solution Imprecise for imprecise exceptions

لا في ال priorities ال software

add. ال opcode ال

ال software

ال software

ال software

Home

→ they both put SCAUSE and put 4 in it

→ Software (if statements) decide from where to

start solving (software priority)

Topic

Parallelism via Instructions

the design of pipelining and parallelism through stages is called ILP

The multi-cycle

5-stage pipeline ← instruction-level parallelism

clocking depending on the time of least-time stage for all clocks

time = $\frac{1}{F}$ and vice versa

The new solution to increase performance:

The Multiple Issue

for example

→ The fetch stage is giving you more instructions

[not only one instruction at a time]

كثافة أعلى من النسخة (تخزيناً) pipeline يوجد كل دورة ع

لكن حقيقةً احنا بنقل replication من ال components وبتساكن كالتالي

فبتساكن بتعمل عليها واحنا بتساكنها بتساكنها بتساكنها

→ * Ideally: in the normal situation, where no stalls, no problems (only sequence of instructions)

in the 5-stage pipe-line: each cycle you finish one inst

$$CPI = 1$$

WB

IF

$$\rightarrow IPC = 1$$

in multiple issue: 2 instructions enters, 2 exits (for example)

$$CPI = \frac{1}{2}$$

دورة 2 وسيفر 2

البرامج
التي هي متفرقة

to avoid using fractions we use new unit called IPC: instruction per cycle, not cycle per instruction (✓ المطلوب)

$$IPC = 2$$

check the example in

slide 65

peak \equiv ideal

Problem in multiple issue & dependencies!!

two instructions have true dependencies in between, are going through stages together, while you can't forward, because they are the same place

increasing the depth of pipeline (more stages / sub stages) is not that best solution

↓
increase the width (multiple issue)

Multiple Issue

① Static Multiple Issue

For simpler computers (few power use)
→ for cores of not high end computing
for things like videos applications

هذا هو العامل المحدد في الأداء، لذلك يجب أن يكون أكثر ديناميكية

↓
This is limiting factor, let it
be more dynamic (by CPU)

② Dynamic Multiple Issue

grouping and all these things, let
them be the job of CPU

→ في الـ instructions يجب أن تكون متجانسة، وأن يكون الـ HW قادرًا على التعامل مع الـ dependencies والـ control signals

hazards، والـ instructions المتجانسة، والـ hardware

But compiler helps a little bit

Static Multiple Issue design

compile \rightarrow instructions \rightarrow packets

Grouping

the mix of instruction types in packets depends on the design, you can't put any instructions together with no restrictions

* \rightarrow Static \rightarrow VLIW \rightarrow 70% من الـ static



فانقره كثير عن static concept
very long \rightarrow تاوانه و اجوده inst و بطولم كانه
inst و صفة ديموم بده و لان تا نفس الوقت



multiple opcodes

① In VLIW designs, instructions are saved in external memory like DRAM or flash memory in embedded system, or in the CPU cache

② we have Execution Units (EU), multiple of them, enough enough to deal with the big number of instructions, and the absence of any one of them causes a structural hazard.

design of slide 67 is containing integer register file
③ we can't replicate register file. ← لا يمكننا أن نكرر

↓
but we have to ~~change~~ modify it some way to be able to deal with the big number of instructions

because:

- we only have 2 read ports - 2 ~~write~~ ^{output} ports
- port for the input coming from WB stage (write port)

But for 1024 bit (32 instruction) at one time, this is problem because the sequential access of instructions on RF

→ increase the number of ports

تعدد المنافذ
العدد الكبير من

RF

↓
in parallel

→ we do padding when no operation can be added. → (giving all value 0)



ضروری لائنیں بالکل صفر

علاوہ اس کے signals ضروری لائنوں

یا 0 یا 1 اور داتا فیلڈ

الیا باقی

We have ² Dual issue, ³ ~~trip~~ tripple issue, ⁴ Quad issue

In slide 70

how we determine the address of first PC of the packet, (PC+X)

it depends on the arch:

- length of the instruction

- number of the instructions per packet

RISC-V with Static Issue



The 5-stage pipeline

one of the most important targets -

- to have enough resources in EX stage to deal with dual issue

Slide 71:

First added unit to the design

EX stage:

we leave the first ALU with its muxes to calculate for ALU instructions & branches.

But, we add additional ALU to calculate for LD/SD addresses.

No structural hazards

Second added unit

add new Imm Gen unit to decode stage

ALU bit 12 of inst div imm

LD/SD inst. ALU inst. extension

For example :-

ADD x1, x2, 45

LD x5, 8(x16)

Third Addition

Duplicate parts of RF in ID stage

الـ ALU 4 أجزاء PC الـ 71 \times 32-bit inst \rightarrow dual issue \rightarrow next PC
الـ 8 bytes \rightarrow 8 bytes \rightarrow next packet \rightarrow 8
الـ 4 \times 8

Hazards with Dual-Issue RISC-V

- ① true data dependency, forward to next packet
- ② local-use dependency, forward to the packet after the next one

\rightarrow wires and units of forwarding should be replicated because we might be intending to forward to an instruction in the same packet, and for an instruction in the next packet.

Scheduling Example // Slide 78

-points to keep in mind

① compiler will take the role of spreading instructions to packets. (static design)

Explanation of code:

x20: the base address of our array
is saved in this register

→ take offset 0 on this address, go to it,
in memory, and load the element stored in that
location

x31: register containing array element.

x21: any scalar

line 4 (decrementing pointer), to update the
target address, so the next time you want to
ld something from memory, you find address
in x20 updated

→ we are updating the address using (-8) shift
because we are dealing with long long
values (8 bytes)

← سؤال واحد عن قيمة shift في address، هو دائما عن 8
byte بالعمود لأن وقت loading بحجم واحد 8 bytes

some analysis

- two instructions → load/store
- 3 instructions → Arithmetic/branch

already imbalanced

- addl inst depends on ldl inst (load use dependency)
 - ↓
 - can't be both in same packet
 - addl can't be in the next direct packet we need (filler) if we can.

- store inst depends on addl inst
 - store can be in the next direct packet
 - Forwarding Unit can solve the problem

- blt depends on addl;
 - each in separate packet.

First scheduling (naive)

| | | | |
|-------|--|-----------------|---------|
| loop: | addi x20, x20, -8 can put Addi here | Ld | cycle ① |
| | ↑ or ↓ can put Addi here | Nop | cycle ② |
| | addi | Nop | cycle ③ |
| | blt | sel x31, 0(x20) | cycle ④ |

addi does not depend on any previous inst, so we can put it on any packet before the packet of ~~and~~ addi inst.

Mistake!! problem in solution

We've changed the sequence of code, the flow have been changed (logical error)

- original correct sequence!

- load element from address x20

- edit element value in x31

- save value x31 in the same current address saved in x20 (the old address)

- change the address in x20 to the next target one

- ~~the~~ new flow of sequence

- we put `addli` before `sd` inst, so that we update the address of `x20` to let `sd` save `x31` value in the next target address not the old address that we were working on through this operation.

اگر: بیٹا دینی کا ٹیپر، انجینئر، اسبقہ علی ال address الی برہا ریج
علی ال sd

new update on code:

`addli x20, x20, -8` → بقسٹہ ال 8 address
`sd x31, 8(x20)` ← دوسرے ال x20

اگر ال sd، دوسرے 8 عسار الی کا ٹیپر ال -8
بدون برہا ریج، لقمہ الی خزینہ ال x20

Hint: if you put the instruction of updating the address earlier, **reflect the effect**! because we basically update the address after we finish using the old address

IPC? $IPC = \frac{\text{number of actual instruction billed in packets}}{\text{cycles used for packets}}$

$$= \frac{5}{4} \approx 1.25 \text{ (peak } \approx 2)$$

min 13

How to solve the problem of low IPC?

why this low IPC?

- filling with Nop
- the dependency (but not as a direct reason)



لماذا سلكوا حيل الورد، حاد رينه
نستخرج inst ناسين نوي فيه

كيف يمكن نستخرج هاي حيل نول
النتيجة؟

التوازي بينهم عن ال depend

Solution: increase the number of instructions!

work on your code to expose more instructions on assembly!



Loop Unrolling compiler's job

Example:-

* the rolled loop

```
for (int i = 0; i < 10; i++) {  
    a[i] = a[i] + x;  
}
```


* the unrolled loop

```
for (int i=0; i<10; i=i+4) {  
    a[i] = a[i] + x;  
    a[i+1] = a[i+1] + x;  
    a[i+2] = a[i+2] + x;  
    a[i+3] = a[i+3] + x;  
}
```

→ This way: reduces loop-control overhead

loop-control: (i < 10 ; i = i + 4)

→ update i

→ compare



* in rolled code: 10 times running code

11 times implementing

control (10 upgrades
11 comparison)

* in unrolled code: 4 times upgrading i

5 " comparing i < 4

← لتوضیح برای unrolling کد الی iterations ، پس عند آجز
تستین ، بهای کامپیوتری که در دستار جایگزین الی ← outOfBound

This operation is done in assembly not C!

Steps of Unrolling in Assembly! Example slide 82!

① Replicate 4 times! copy/paste

```
lcl x31, 0(x20)
addl x31, ---
scl x31, ---
```

} → $a[i] = a[i] + x;$

```
addli x20, ---
bit x22, ---
```

} $i++$
 $i < 10$

→ this part is not needed in the first 3 replicas

we only want to change i value after 4th copy, and to compare there!

② Remove unneeded loop overhead!

Still having some problems

4 copies are still identical instructions!

they all do $a[i] = a[i] + x;$

while we want them to work on 4 different

elements of array ~~that $a[i]$ precedes!~~

that $a[i]$ ~~precedes!~~ play on offset!

③ Modify instructions

* modify offsets in this case!

-8, -16, -24, -32

problem!! we are only working on x31!!

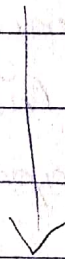
-logically, no problem.

- in case of scheduling, all your

packets will depend on x31, which causes dependencies problems

التتابع على
Packets
متتالي
sequential

ناتج سبب
dependencies



④ Rename registers! (complete renaming)

- give each group different register name

متوحد اللى جيتي اعمل renaming؟ هل خبير اعمل 64 unrollings!

واحد على واحد registers محدود لعدد الجانج خاصة، ففي على limit

على renaming من برادق! قلنا قبل على جيس، الوبيلار على بطور
dependencies، بيتر على data path 8 inst مع جيس، ولا

2 inst مع جيس
← 2 inst ← بتر اعداد كثيرة، وول مايلو wider جيس اصب

⑤ Scheduling instructions:

example continuing slide 40:

* take into account these things:

- we have local use dependency in each

group between \rightarrow ldl \rightarrow addl

filler packets \rightarrow addl since ldl is a filler

- true data dependency between addl, scl in each group (no need for a filler, just put each one in a packet)

dummy solution:

| | | |
|------------------------------|----------|---------|
| | ldl x28 | cycle ① |
| | ldl x29 | cycle ② |
| ال scl لا يتم تأجيله | addl x28 | " ③ |
| يجب ان لا يملأ | addl x29 | " ④ |
| في حزمة تلوها | scl x28 | " ⑤ |
| directly in the next packet. | scl x29 | " ⑥ |
| | scl x30 | " ⑦ |
| - we can change this | ldl x31 | " ⑧ |
| to win some cycles | addl x20 | " ⑨ |
| | addl x31 | " ⑩ |
| | blt x22 | scl x31 |

Optimized solution

cycle (1)

" (2)

" (3)

" (4)

solved in slide 91 ←

" (5)

" (6)

Buf!!

" (7)

a problem

" (8)

← لاحظ أننا اننا اعددنا الى اول packet

" (9)

ما حسبنا حسابنا اننا هي بتقول على x20

" (10)

التي كل ال lca و sdc بتقولوا منها

فلاحظنا اننا ← offset نحتاجه ل lca و sdc

نحسبنا طريق اذا كاننا نحتاج الانتقال من ال lca

نلاحظ أننا نحتاجه ل lca و sdc

التي نحتاجه (-32) يعني ان offset ل lca

x28 lca و اول packet

✓ packet نحتاجه نحتاجه

lca و نحتاجه ال

$$-32 + y = -8 \rightarrow y = 24$$

inst: lca x24

$$-32 + y = -16 \rightarrow y = 16$$

inst: lca x30

$$-32 + y = -24 \rightarrow y = 8$$

" lca x31

$$-32 + y = 0 \rightarrow y = 32$$

" sdc x28

$$-32 + y = -8 \rightarrow y = 24$$

s sdc x24

$$-32 + y = -16 \rightarrow y = 16$$

= sdc x30

$$-32 + y = -24 \rightarrow y = 8$$

s sdc x31

Static scheduling ends here

* Back Deeper in Dynamic Multiple Issue

"superscalar" : more than one instruction entering at one time (this is done by scaling the number of instructions entering a stage / we are scaling it.

→ Avoids the need for compiler scheduling, but it may still help us in operations like: code loop unrolling and exposing parallelism.

code semantics ensured by CPU



الاستيعاب الآلي

إدارة نوافذ الـ

offsets

السcheduling الآلي

إدارة الـ

نوافذ الـ

logically

this is the job of CPU now

To understand marketing sheet, Intel
ARM, AMD, RISC-V, MIPS web sites

(or words like (superscalar, out of order,
in order) processors will be used

superscalar: Dynamic issue, multiple instructions
being fetched, and handling + managing them is
CPU's job

out of order: our CPU, can do execution for all
instructions coming from IF, ID stages, in orders that
never relate to the original order these instructions
came within

- Out of order execution } we all return to something
- in order Fetch } called (In order commit)

out of order expression only relates to
execution stage ✓ Some thing have to be
done to ensure logicality after
this stage !!!

what we mean by out of order?!

سلاسل (مثال) : دبی اسج ال instructions کا ہرے کسے

یہا تھر بنا ال instructions

can start sub while add is waiting for lcl

logical results جس دھیت اداوڑ عال

دھیت آخر، سو تھر!

- با طائے ال processor بنا ال sub وال and قبل ال
lcl وال and local-use dependency

دھول ماہیم dependency

- عتہ ال sub دھول ماہیم lcl وال and عتہ ال dependency
ہیم

اگر دھول فرتہ سیکڑے ال اسٹال ال

out of order

Dynamically Scheduled CPU (design style 94)

→ As major stages, major view on design, all processors of this type have something similar to it, not forgetting that any stage in any design can be divided to many sub-stages.

→ (1) Stage (1): **ID + IF** stages are combined in one stage (Instruction fetch and decode unit) → control unit

→ loading instructions

↓
depending on each inst type we'll spread instructions each to the suitable unit (available execution units)

(2) Stage (2): **EX**
functional units ↓

number of units

and their types are different from design to the other

processor

Stage: Issue

Reservation Stations

- issuing instructions ~~are~~ is done here
- like a buffer (they are registers) to let instructions wait until the units of executions become free to work on new inst, to avoid structural hazards of not enough resources.

- even if EX units are free, we may leave inst. in reservation units to wait until some problems of dependencies are solved.

← reservation stations
← RF الوجود ال stage IF/ID

↓
if the inst. came and there is no dependency, it reads registers operands and puts them on the beginning of available units

← for this instruction is ready
← EX units

scenario: Reservation station

نقرض ان يكون لدينا 4 inst على RS الاولى ، و 4 inst على RS الثانية ، والثانية خاصة ال EX توجدنا وهنالك قطعة ال integer unit ، بالعبارة ال RS الاولى رسولنا هو دفتر النقل ال inst ، بالعبارة ال RS الاولى او ديها عال الثانية التي هي وحدة ال RS 1 + RS 2 : two RS in the left of slide 94, opposite to integer units.

اذا بناتنا design فيه هياي اللامانية ، بنصاح RS وحدة يتفرق ال RSs independent ، حسن هاد طاني مني موجود ارمال برانين تبعنا RSs independent ، وكل وحدة خاصة بال EX unit بالعبارة ال

كافة النتائج، انه ال results ال CDB ترتيبها من اسي
مع ترتيب ال inst ال كما هو !

↓
call inst enter
↓
ready inst start to execute

↓
CPU is controlling and choosing
↓
no guarantees on maintaining the
original order

Because of this, at the beginning of commit
unit there is a stage called ROB

* ROB stages re-order buffer → بترتيب ال ترتيب

↓
السابق لترتيب ال inst ال ال
How the ROB know the correct sequence?!

↓
هو بيقتر نتائج ، سو يد يفهم هاي نتائج سو هو

وليف يدوعرف ترتيبها!

بعيداً عن ال ROB احنا بيقتر ترتيب ال inst من

ال PC وال PC بيجيب ال inst من ال inst memory

الاي فيها ترتيب ال inst بالترتيب ال ال ال ال ال

1 1

→ Reservation stations + other circuits do the job of tracking the right sequence of inst. depending on their PC's

Single Issue Example slide 46



multiple issue.

Code Analysis:

- local-use dependency → ld
- add
- true data dependency → sub
- andi

Single Issue? only one instruction

- EX stage sometimes limits the number of instructions

EX stage sometimes limits the number of instructions
Queue
2 instructions

min 41 ✓ lec 13

Wares Recovery

out of order : describes EX stage

التي لا يتم تنفيذها في الترتيب الذي تم كتابتها فيه

Committed by order

C stage is sequentiated
ان C لك inst التي
منوع تكون قبل ان C ان inst
التي قبلها - يا اما بنفس ال
يا اما جالي بعدها

Register

Register renaming is replaced Slide 47

Dynamic reservation stations

ROB reservation stations

التي هي في شكل RS و inst و RS و حجز بيوتها

التي هي في شكل RF او CDB او
Reorder Buffer (ROB)

باللحظة التي ياحد منها القيمة في احدى هذه

Registers التي كانت تستعمل في

التي كانت في ال cache ، و ان كانت في ال cache

التي كانت في ال registers

dynamic reservation stations registers

renaming

explicit في ال compiler ، و ضمناً موجود في ال HW

Recovery :-

RAW (true data dependency)

this dependency can exist for any type of design.

But the WAR (Name dependence)
WAW (Output dependence)

↓
they may only

happen in

Multiple issue

out-of-order

pipelines

* Name dependence: WAR

should: Write After Read.

Code analysis slide 98

- true data dependency \rightarrow mul \rightarrow addl.

waiting for x1 to get ready.

* lcl does not depend on addl, but it will write on x5 before addl reads the correct value of x5 \rightarrow in 5-stage pipeline this ~~was~~ wasn't a problem

In dynamic processor (Superscalar) deg 3

① mul and lcl are ready to move from RS to EX units.

② ~~mul~~ mul is a long instruction, it will take time to finish and find x1 which is the second operand of (addl), hence, addl will stay more time in RS

③ through this time lcl will write on x5 before addl ~~is~~ reads the correct value of it.

data dependency \rightarrow $\{ \text{mul, addl} \}$ ✓

This problem will be solved automatically ~~is~~ by the in order commit.

- in order fetch
- out of order EX
- in order commit.



* Output dependence (WAW)

الاعتمادية في الذاكرة و في الذاكرة
 بالترتيب الطبيعي كما في مشكلة في الذاكرة والذاكرة

→ Because of multiple issue and out of order

① l1 has finished execution earlier than mul, then the latest value of mul will be from mul not from l1, and this is wrong logically

It will be solved by in order C

remember: W stage in Dynamic is the stage where results are written on CDB

The In order Commit + Reorder Buffer can solve these problems.

Speculation (Guessing)

→ Branch ✓

→ Local ✓

→ Remember, prediction implementation is done in F stage

→ resolving prediction: depends on design

5-stage pipeline: ID ← - it may be in M

ID M WB WB WB ← - " " " " EX

- " " " " ID

checking if prediction is correct.

In superscalar: resolving branch is done in EX stage, because there we implement the subtraction between the two operands of branch inst.

→ In 5-stage pipeline, flushing instructions was easy, but in multi-issue, it is not.

What's the solution?

Let the control unit in EX stage of branch to give (nullify), to do not allow commit stage to work for fall-through, instead of C, → nullify

signals of next stages set to 0.

In other words:

- MEMRead control signal: 0
- RegWrite " = : 0

after Nullify:

- fetch new instructions directly after the cycle we determine that the prediction is correct or not, not after nullifying.

local speculation ~~local~~ problem (Pointer Aliasing)

it happens most times in programs full of pointers. ✓

example:- slide 106

- 3 inst are going through F, I...

(dec 3) → lcl
→ sel
→ lcl.

- first lcl is done peacefully
- ~~the~~ sel is waiting for x1 to get ready from first lcl, but second lcl is not waiting for anything.

second.
- sel writes in RS, while lcl is about to finish execution

But lcl will read from the same address the sel should write a new value in [lcl should use the updated value, but this not possible in this situation]

Speculation and Exceptions

Exception Example slide 111

what is the exception of C?

Remember: we may have these exceptions

F stage: out of range address (when loading instructions)

D stage: invalid opcode

E stage: overflow

C++ → it saves on RF and memory
then exceptions are from type: out of range access.

after detecting the exception in cycle 7

- soonest time we can move to exception handling code is cycle 8 (perfected design)

- first couple of instructions to fetched ^{exception handling code} are gonna be of type (sel) because we should store our values of previous instructions to keep track on our program. (context switching procedure when moving from code to another, save previous registers)

Why Do Dynamic Scheduling? slide 112

→ Remember: we have a lot of ISAs

↓
with different latencies
+ Hazards.

ISA (1): x32 or x32-64

↳ Intel / AMD } Both can run windows

+ Both can run same programs

↓
because of same ISA

ISA (2): ARM V7 (ARM company)

ISA (3): RISC-V

ISA (4): MIPS

ISA (5): PIC

- Why not just let the compiler schedule cache?

Answer:

- take into account that even one company can produce more than one generation and design based on one ISA.

Intel: - Pentium

- Core

- Core i3/i5/i7/i9

- XEON

- Silver / Bronze / Gold

To let the compiler schedule instructions into Issue packets, you should provide it with enough knowledge about all these designs. unrolling و scheduling

كتابة كود بايزر بهاي (ترياعه هوس) دوسن افرا

- Not all stalls (caused by hazards) are predicable by compiler (can be handled by compiler)

cache misses can't be handled by compiler because it does not know what exists in registers which is not helping to handle this situation.

1 / 1

- for modern CPUs (cache included in the design) : High-end designs.

- instructions are in instruction cache
- data is in data cache

- Not modern CPU's

- main memory (DRAM or Flash)

↓
always the memory

closer to CPU

- Can't always schedule around branches.

- compiler does know if the branch is taken or not, this is detected at run time only.

- Delays, limited bandwidth of the memories (memory delays), can't fetch fully parallel instructions together, all functional units are waiting for ready instructions to be executed, it is hard to keep your pipelines full all the time.

We are better than before, but not as wanted to be.

Chapter 5: Memory

Computer System: - CPU and

- Buses

- Memory subsystem

- We assumed in previous lectures, M stage takes only one cycle, which is not realistic

Assume that:

our design is about: - CPU

- RAM only (no caches)

↓
M stage will take from tens to hundreds of cycles.

→ Our DRAM is a middle way to transport data to CPU instead of slow Hard disk where M could take millions of cycles.

→ CPU speed = 4 GHz (assumption)

→ hard disk at highest case: 7200 round per second

Huge difference!!!

↓
Mechanical

But while using DRAM, it is still bad with 100's of cycles of M → Solution: cache between CPU, DRAM
holds data

في سرعة الذاكرة cache انه خاص عليها اهم دانا بالتحديد للانتم في
تقليل زوايا وخطى على ال RAM

تأكيد الا فزوع وزيج صيد DRAM بس بدنا نقل ال Frequency

Principle of Locality

هو (مفهوم) local

Temporal
↓
time

Spatial
↓
space

الا جاد اللى بيبدأ كل وقت
time + space

Example :-

```
int x = 0;  
for (i = 0; i < 1000; i++) {  
    x++;  
    a[i] = a[i] + 5;  
}
```

سؤال : سوال variables اللى كل لغة قاعد ال loop جيتقل عليهم ؟

i و x

اخاف عن جيتقل عليهم كل لغة ، وكل لغة ولغة بديهم فولدق ، فنية
وكلهم !
→ each one of them has temporal

locality (efficiency in time)

رصيداً عن فنية بديهم كل مرة جيتقل عن نفس ال variables

كيف يترتب ال arrays في ال memory

① row major ← Java, C

② column major ← Fortran, Matlab



تقسيم الذاكرة

في الذاكرة ال memory hierarchy ال hidden hardware
البيانات ال data ال hardware

In slide 5:

local storage

(Flash)

→ like USB sticks



hard disk

↳ hard disk

local storage (Magnetic) HDD

But Flash is SSD

Question: if we invented a technology that non-volatile memory, cheap, and can response to all local and store instructions as much as they need responses within one nano second

it will replace the HDD or SSD, and caches, DRAM and we remove them to put more cores.

But RE won't be replaced (it is the core of CPU)

Example

If 100 memory locations are loaded to cache,
But the CPU could find 60 of them

$$\text{Hit ratio (rate)} = \frac{60}{100}$$

$$\text{Miss ratio (rate)} = \frac{40}{100}$$

what is the meaning of hit and miss?

we loaded [it], because of spatial locality,

a block of 16 elements of array a is loaded,

when the CPU tried to access location 17 of

this array, it could not hit it in cache (miss occurred),

so, it went to DRAM and loaded another block of

16 elements, ~~now~~ where location 17 is the first

element of it and this takes tens of cycles.

Section III Memory technologies

SRAM technology

← ال transistor هو عبارة عن switch يتحكم في مرور التيار بين القطبين
بواسطة الجهد المطبق عليه



عند تطبيق الجهد، يمر التيار، وإلا فإنه يحجب التيار. [it passes the current or blocks it

بال voltage يتحكم في مرور التيار.

ال 0 و ال 1

block كإشارة

- Fast but not dense

1×10^6 transistors \rightarrow 6 for each bit

$$\frac{1 \times 10^6}{6} = \text{about } 166 \times 10^3 \text{ bit}$$

↓
not dense

in embedded systems we put the instructions in flash memory, and data on SRAM memory \rightarrow you can find chipsets of SRAM technology ^{which} you can interface them to CPU.

on the chipset of slide 9.

Example: our SRAM size is 64K bit, how many chips you need? $\frac{64K}{16K}$ chips. 4 chipsets

the \overline{CS} is connected to a decoder which decodes the address, and depending on that we choose which chip set to activate.

- output of decoder is a binary number where only 1 bit is 1, and the rest are zeros

In five-stage pipeline, in stage M: there was a control signal named Memwrite, this signal is connected to \overline{WE} in SRAM

DRAM technology

it saves data on capacitors not transistors to save data permanently

← مستوى الذاكرة ، مستوى الذاكرة ، مستوى الذاكرة ، مستوى الذاكرة ، مستوى الذاكرة

← مستوى الذاكرة ، مستوى الذاكرة

Problem (electrical)

the capacitor will discharge overtime, so it is not able to return our information

the solution is to read all saved data again (capacitors) periodically, because this lets capacitor charge again so it keeps saved value. (refresh cycles)

The refreshment makes the controller complicated, now we need more circuits to deal this operation of refreshment now the controller reads, writes, & refreshes

A problem on refreshment.

conflict between refresh cycle and read instruction on DRAM. (~~write~~ they are going to happen in the same moment.)

→ Designs have different points on giving the higher priority to one of them to be done first.

DRAM access sequence:

→ address of load or store instruction is analyzed by the controller inside the DRAM

→ after analyzing the address it is gonna be divided into two parts, a part to identify row, and one to identify column.

example:-

ld xt, 8(x2)

address in memory = ~~8(x2)~~ $8 \times 2 + (x2)$

this address is put on bus to reach controller

→ address parts are saved in row address latch and column address latch, to save them from being lost because of the next coming addresses

Big image: bus, memory, CPU ✓



← ال CPU ال address

← ال cache ال address ال hit

← ال address ال miss



← ال controller ال starts working
← ال controller ال sends a request to memory to get the target address



being put on bus

it reaches the controller



← ال address lines
64 lines or 32 lines

32 bit based design is not enough to address
install more than 2^{32} / 4GB RAM on computer
which is not enough, so we moved to 64 bits
based design

← ال address ال 64 bits ال physically
48 lines ال signals ال

Now, the controller divides the address into two parts (row, column)

Basically, why we ~~we~~ divide the address?

Because the data is stored in memory as rows and columns (tables)

4 Back to RAS and CAS

→ after the controller divides the address, (row, column)

now we want to save them

→ give a signal named (RAS), to save (row address in latch) then pass it to a decoder, to identify the row

(this operation takes time (TRAS: time to assert RAS (to latch and decode))

when this ~~time~~ stops:

→ CAS starts (TCAS)



column addr decoder

column addr decoder

decoder



then the data is put on data bus to go to CPU

Note: the spacing between data 1, data 2... is now shorter than the previous design

(higher bandwidth) → improvement

low bandwidth → classic DRAM

higher bandwidth [get a row, put in row buffer taking elements from it - column by column]

→ fast page mode DRAM

↓
FPM DRAM ??

Fast page mode DRAM offers improved speed over standard DRAM. When memory access is performed within the same address row (page), a single precharge is required for only the first access. Subsequent access within the page can then be made without the time penalty of additional precharge cycles.

Although it has lost popularity in some areas to faster technologies such as double data rate SDRAM, But FPM DRAM is still a very cost-effective solution when the highest performance is not required

- A problem in FPM DRAM ?!

knowing that a row is entirely coming because of special locality

But! what is the thing that control the values of col.1, col.2, col.3... (addresses of columns) in slide 13#

تلك القيم التي هي في الذاكرة (address) في
الجزء الذي هو decreasing من col.1 - col.2
في manual work من جهة أخرى!!

بما أننا نعلم أن كل شيء في الذاكرة يجب أن يكون في مكانه
بمجرد أن نكتب column address وهو يجب أن يكون!

بطل!!

لأنه يجب أن يكون لدينا في الذاكرة أول عود و نحتاجه عن counter و نحتاجه
عن الذاكرة counter ← clock و الذاكرة يجب أن تكون automatically

و يجب أن نضع ال elements في الذاكرة (row buffer) في الذاكرة
في الذاكرة (row buffer) في الذاكرة (row buffer) في الذاكرة
التي هي في الذاكرة (row buffer) في الذاكرة (row buffer) في الذاكرة

مما زاد عليه الزود اسرع!!

→ Synchronous DRAM: with clock

→ Asynchronous: without clock

Recovery for the name

→ RAM: you can load data from any where in it (Random Access Memory)

→ DRAM (Dynamic): in the design we have 1 transistor with capacitor, and it has refresh cycles (a lot of dynamic work)

→ SDRAM (Synchronous Dynamic): clocking for the counter that gives higher bandwidth and speed advantage.

DDR: depends on the concept of reading two columns in one cycle, one on the positive edge, and one on the negative

↳ while in SDRAM, we only were able to read one column in each cycle, either on the falling or rising edge.

Nowadays we have DDR1, DDR2, DDR3, DDR4--
[differences between them depends on the speed of clock, starting from 66 MHz or 200 MHz, to 4200 MHz or 4660 MHz speeds [solving the problems of electromagnetic interference and such problems, to be able to reach these velocities of RAM]]

Higher speed \rightarrow higher frequency \rightarrow reading
Higher Bandwidth \leftarrow doubled data.

There is also GDR \rightarrow GDR SRAM
 \rightarrow GDR DRAM } Not that popular

↳ Super computers application.

Main Difference between DDR and QDR :-

QDR uses two clocks, with phase shift between them

- 2 clocks (positive, negative) for each
- two ports
- then you can do 4 operations

But!!

- the first clk is dedicated for reading data

- the second clk is " " " " writing data

→ [in DDR, data had only one direction, you either write or read]

- Reading/writing can be done at the same time in QDR (dual ported)

↓ ↓
read write
port port

→ Maximum Usage of QDR only happens when you are reading and writing in the same time, throughout all your work.

Quick Recovery

There will be

□ controller

to receive the addresses of CPU and cache sub-system requests



This address is transferred to be divided to two parts ^{rows} _{columns}

in slide 17:

→ Row address latch and decoder

→ each decoder is 14 → 16384 decoder



14 decoder

16384 outputs
only one should equal 1 after decoding.

→ column decoder

10 → 1024 decoder (only 1 of these 1024 lines should equal 1 after decoding)
this means that each row has 1024 elements

After putting data on buses (on parallel buses), then it reaches the north bridge, and the controller in the motherboard (in modern PCs and laptops) from where it goes to the CPU and be put in the cache sub-system

-DRAM Generations

-capacities are doubling, because DRAM are ~~more~~ similar to CPUs, where the design is based on resistors, and we know that resistors are following the principle of Moore (shrinking the size of resistors, doubling the capacities each 1 year and half)

it is pretty easy to find ^{Byte} 32GB capacities these days on sticks

What is the important point now?!

-improving TRAS, and TCAS ✓

Reminders

from the moment, Row Address Strobe is asserted, and address is analyzed, latched then decoded, then the row is moved from row array to row buffer → this is TRAS

→ Asserting CAS, latching, and decoding column then putting data on bus from the row, this is TCAS.

Why the improvement on TCAS is better than TRAS?

↳ moving from classic DRAM

↳ to FPM DRAM

↳ to SDRAM

have improved the work on columns

↳ and this is improved more when we moved to DDR (reading two columns / cycle)

In slide 20#

talking about DRAM Banking

→ from the beginning, why Banking?!

→ it increases the bandwidth, and this can be proved by numbers

But how Banks are used in design?!

→ Back to slide 17#

reaching the procedure again: ↓

① row address is found

② row address is now an input for the decoder

③ after getting the output from the decoder, you can read more than one row, each is from

different region (Bank) in the same time (in parallel)

→ this makes the control easier and helps to distribute data efficiently

But how to decide which row of those has to be read?

the Bank control logic (decoder) decides which row of these has to be read.

احتمالاً، توزیع ال rows علی اکثر Bank و ال banks
البرانی، کلیم بین ← 16384 rows، اذا ایا توزیع ال bank
وینا row، 1000، ال Bank و توزیع ال row ال
کریه ال



So you have two controllers, decoder to
identify row location, the other decoder
is to identify from which bank you should
read that location

احتمالاً، فرضاً، rows على أكثر من Bank، وكل ال banks
البراني، كدم، 16384 rows ← إذا اختلفت ال bank
و.ب.نا row، 1000، ال Bank، ال row ال
كريبه ألف



So you have two controllers, decoder to
identify row location, the other decoder
is to identify from which bank you should
read that location

studying slide 21

Figure (c)

scenario on the communication:

Assume the processor is doing ld/seq , and is heading to cache, to miss or hit the target address (Hit: data is loaded from the cache to RF and the processor works on it)

But if miss?!



the cache sends a request to the higher level memory (DRAM in laptops/mobiles/computer systems), the wanted address is sent on the bus (requesting from memory)



and the process of RAS, CAS, decoding that are done on the addresses in the memory are understood before.

then the memory returns the ^{requested} data on the bus to the cache sub-system, to let the CPU work on it

So what should be found here?! (calculated)

Bandwidth and miss penalty

Miss penalty: in the case of trying to access some address in the cache by the processor, but the wanted ~~data~~ content is missed, you should go through the bus to memory (hence the time ~~taking~~ elapsed to pass the address from the cache, and doing all needed processes in the memory to return the data to cache)



This is the penalty

Bandwidth: the amount of Data sent on bus the time needed for this data to pass the bus.

The example in slide 21: miss needs to reload

16-bytes block from the memory

→ Remembering (words vs bytes)

standard (most used) is word = 32 bit = 4 bytes

intending to read 4 words from memory

But! what are the specifications of memory?

قائمة
نوع
integers
قاعدة
نوع

in figure (a)

memory is said to be: one-word-wide memory organization

↓
what is that?!

width of each location is one word (4 bytes) 4 bytes by 4 bytes

↓
locations are internally

↓
which means in this example that our 16-bytes are divided into 4 locations

* Bus has the width as same as the memory (one word wide)

↓

In summary:

- 4 bytes bus / 32-bits; means that the bus can hold the entire word while transferring it

(there is no any bottleneck here)!!

↓
the bus can transfer the entire word in parallel

- 4 bytes wide memory dedicated.

- classic DRAM → needs request to load each ~~byte~~ word from memory

Example applied on figure (a):

- 1 cycle needed to request a word from memory
- 1 cycle to send data on bus from memory to cache
- 15 cycles needed to process data and get it from memory

$$\text{miss penalty to get 4 words} = 4(1+15+1) = 68 \text{ bus cycles}$$

single request for each word

4 is multiplied because it is a classic DRAM where dedicated complete work is needed for each word.

we call it (time delay)

$$\text{Bandwidth} = \frac{\text{number of bytes}}{\text{miss penalty}} = \frac{16 \text{ bytes}}{68 \text{ cycles}} \approx 0.24 \text{ B/cycle}$$

Target!! improve the previous design?!

How?!

Suggested improvement: to increase the out
width of memory, to be as wide as the
number of words you usually need to read.

↓
usually need to read 16 byte: let it be
16 byte wide (as an example)

→ bus should be widened also

→ even in classic DRAM, you still need one
request for each 4 words (16 bytes)

that's the point!! (one address
needed)

no redundancy ← to load the
the operation is done once block of 16 bytes
one request to get the 16 bytes (4 words)

→ and then the cache (using a multiplexor) chooses
the bytes that CPU needs ~~specifically~~ specifically

multiplexor

→ multiplexor slide 22 Figure (b)

But! What is the problem of this design?!

Quick Recovery

→ increased the width of bus

the bus is the electrical real lines that transfer data from the memory to CPU and cache subsystem

physical lines
mother board ← die

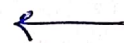
~~what~~

→ because of this widening, now we use 128 physical and electrical line, instead of 32

↓
in routing (how you

put them on mother board), will take a wide space

80:



① we'll increase the size of motherboard which means higher cost.

② electro magnetic interference (EMI)

which causes errors in the transferred data.

↑ this is the problem!!

what to do now?!

→ now, saving time in row decoding, that takes the longest time in the whole operation, only it takes time at decoding the row address once for all banks

Now:

- ① send the request in 1 cycle (send the address)
- ② in 15 cycles, you can get all the rows containing the target data from the same row from all banks

now the problem is (How to return data!!)

we can't put all the 16 bytes once on the bus (it is 1 word wide), and we should transfer 4 words from each bank.

So we use a multiplexer, and in each turn we choose one bank, as we are saying:

bank 0, put your data on bus and transfer (1 cycle)

" 1, " " " " " " " (1 cycle)

" 2, " " " " " " " (1 cycle)

till bank 3, hence we have 4 cycles

③ 4 cycles to transfer data from banks to cache

↓

20 cycles miss penalty!

→ wide, bandwidth
close to the ideal
without problems

↑ this is called bank interleaving

Can we do better? yes!

the previous banks system is not DDR, it is just SDRAM, when it is DDR (where you can send two words at a time) at a cycle)

now we multiply each bank by 0.5 not 1 because in one cycle we transfer the data of two banks

$\rightarrow 0.5 \text{ of } B_0 + 0.5 \text{ of } B_1 + 0.5 \text{ of } B_2 + 0.5 \text{ of } B_3$

↑ positive edge ↑ negative edge ↑ positive edge ↑ negative edge

↓ cycle ↓ cycle ↓ cycle ↓ cycle

= 2 cycles only

slide 23 ↑

↓
✓ one on positive edge
✓ one on negative edge

→ if we used the QDR, we multiply by 0.25 not 0.5, because we transfer 4 words per cycle

But!! this calculation success when we are transferring data in both directions, (in and out of cache), because QDR transfer 4 words only in both directions. It can't transfer more than 2 in one cycle per one direction (read or write direction)

So it stays 0.5 not 0.25 because QDR is DDR for each part (write or read part)

Quick Recovery;

SRAM (of cache and registerfile), and DRAM

→ but these types are nonvolatile, but ~~we~~ we need a permanent storage (so we have the Hard magnetic disks) and the flash memory

→ For flash storage; we don't need battery inside the SD card or the USB flash to retain frequently to keep the data inside it

→ even when applying it to system it does not consume huge amount of energy

→ the numbers written on chips of flash of memory are the speed grade

the speed grade of flash should be suitable for the application

(like HD and 4k cameras, that handle huge amount of data, they need high speed grade to be able to handle all video frames and save them)

→ the low speed grade can be used for numerous of images only

Flash types

→ the modern SSDs, there are algorithms written on the controller hardware inside the SSD, which does some thing called wear leveling → \rightarrow \rightarrow

↳ We can't keep writing on the same location and clear, write and clear while we have empty capacity, ~~when~~ which makes this location wearred while ~~the~~ other locations are still good.

if we don't use this algorithm we loose the capacity of half of our harddisk

~~we~~ wear leveling makes it live longer (SSD)

if removes data to less used blocks

HHD : hybrid hard drive

HDD : Hard disk drive.

Disk Storage

The parts & interfaces

→ زجانه كانه في اسبي اسه ال SCSI - احد انواع التوسيلات التي يتوصل ال HDD بال mother board

وقبل فترة كانه ال ATA هو الذي كان يستخدم عن parallel port (هوائي 40 pins) هاد هو الذي كان يستخدم عن SATA.

→ يعتبر هو ال bus اعداد ال interface - لاسلكي، لاسلكي، ال standard في نقل اجهزة التوسيلات ال mother board التي يتصل ال internal buses

The embedded system

→ يتكامل بكل الاجزاء على انظمة

→ algorithms خاصة ال يعرف ال system ال

address التي كانه ال computer ال location ال cylinder (منه address ارقام binary بطور معين ال levels)

cylinders

→ physically made from ferromagnetic material → مواد دالم منها

التي تتكون منها ال مواد

→ the iron inside this material controls the magnetic directions (the iron has magnetic characteristics), so if we directed iron atoms to north, then we have logic 1, if directed to south then it's logic 0 → such thing happen

* we have more than one layer of cylinder to increase density and capacity. (they could pack more stuff more)

→ so it is easy these days to find 8TB HDD

our data will be on the surfaces of these cylinders

sector sizes:

the ~~sector~~ sector does not have a constant size (capacity) it can change

→ you change it from the laptop from

the allocation unit size

the default is usually 4096 bytes (4k ~~of~~ per sector)

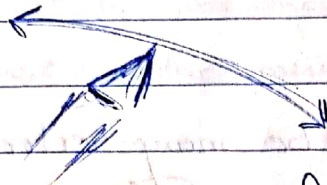
* ideally we hope that the sectors that holds the data of a video for example are existing directly next to each other, But this is not happening actually.

The parallel heads

→ each can move (right-left) on one level

two directions

from the smallest cycle to the biggest one



to reach the rest of sectors, cylinders (disks) will move around so that the head can read all sectors

The speeds of disks moving around are main 3:

5400

5,4 لاطا نسيًا من اللابوبات

لانهم حاسب حلالو طاعة

لانهم motor ويد حاسة عاية

ويهلل الجارية

7200

standard

for all

PCs

above

11000

or 15000

for special usage
(like servers)

2) لاهتر ارج، ليكائنية ل 5400

التي حدة ان 7200، احسنه ل

اللابوبات التي يدرك كثير حتر

زي desktop

فيقتل احسنه لية تلف

* The thing that controls the rotation and the location of each head is the controller in the embedded system

Disk Sectors and Access

Size of sectors (tricky)!

→ if it was so small, there will be many sectors which makes the management harder.

(کئی سیکٹرز کے ساتھ منیجمنٹ سہولے)

→ if it was so large, there will be wasting of capacity some times, because if these times you have one file that is all about $\frac{1}{8}$ the size of the sector, you will waste the rest of available capacity this way

(بے کار سیکٹرز کے ساتھ)

Trade off

→ so then the default size is 4096 (4 kB), as the better size

→ like the EDC (the parity / error detection code)
ECC:

→ errors may come from ~~electrical interference~~ (flipping bits)
→ we have to solve these problems (detect them)

Because of the magnetic interference *

The problem about EDC (the parity: odd and even)
is that it can detect odd number of bits only

→ if 1 flipped bits it can detect it

→ " 3 " " " " " "

→ " 5 " " " " " "

it can't detect 2, 4, 6 ...

* But we want to detect all problems and errors
and hide them, so we found ECC

* there are gaps between sectors to help the HDD recognize the beginning of a sector and its end.

Scenarios

We want to watch the movie, so we open it, and the movie player loads the frames of the movie to display.

→ at the beginning "miss" will occur

→ it goes to cache: miss also

→ " " s DRAM: miss also (called page fault)
↑ something like miss

→ it goes to HDD

→ it sends the address

(if the hard disk was idle, and not busy with another application requests)

Hopefully, the HDD is idle, // there is no queuing delay

- either we have queuing delay (time or not)

we will have these three things when starting to response to request.

- ① seek time
- ② Rotational latency
- ③ Data transfer
- ④ Controller overhead.

controller overhead: time needed to

→ analyze the address

controller of
embedded
system.

←

→ determines its location

→ " " " " the angle of location

→ " " " " the track and sectors

↓ to turn it from binary address to physical locations on hard disk

→ operation of moving the head. (find track)

Seek: time needed to move until you reach the target track

الوقت الذي يحتاجه الرأس للوصول إلى القرص المطلوب، average seek time

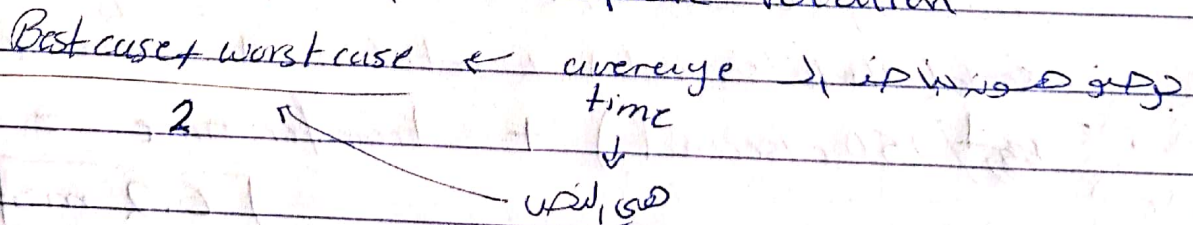
الوقت الذي يحتاجه الرأس للوصول إلى القرص المطلوب

→ operation of rotating the disks (find sector)

Rotational latency:

Best case: accidentally the head is pointing at the wanted sector

Worst case: we need a complete rotation



now recalling the sectors depends on:

- ① data and sectors length
 - ② data transfer rate
- (سرعة قراءة البيانات) (مقدار البيانات)

Disk Access Example.

512 byte : ~~data length~~ sector size

15000 r/m \rightarrow 250 r/second : ~~rotation~~ rotational latency

4×10^{-3} seconds : average seek time

100×1024^2 Byte / second : transfer rate

0.2×10^{-3} seconds : ~~controller~~ controller overhead

no queuing delay : idle disk.

\rightarrow Average read:

queuing (0 seconds) + controller (0.2×10^{-3} seconds)

seek (4×10^{-3} seconds) + rotational average ($\frac{1}{2} \times$

$\frac{1}{250}$ r/second) + transfer time =

6.2 ms

مثال

ثانية لكل دورة

ثانية

ثانية

ثانية

How to extract the transfer time for 512 Byte sector, depending on the previous information:

100 MB \rightarrow 1 seconds to read

512 B \rightarrow ? seconds (call it x)

$$100 \text{ MB} * x = 512 \text{ B}$$

$$104857600 \text{ B} * x = 512$$

$$x = 4.8828 \times 10^{-6}$$

تقريباً، أي 0.005ms

أو الميكروثانية

Note:

إذا كنا نقرأ أكثر من sector، ونسبة "تأخر" كل sectors لك
دراهمنا، فإننا نحتاج حساب الوقت بين تفرغ كل واحد
لنا حسب كل data بدل 512 البايت
transfer

لكننا نحتاج برصبة مرة واحدة.

لكن لو افترضنا ان sectors هي الى او في اي وقت ما تكون، و
بسبب عدم كفاءة كل واحد، فكل واحد sectors فكل واحد
وال controllers التي تكون عليه وفي بعض الأحيان sectors
الخاصة بالآلة.

to keep track how the data is distributed.

seek time في القرص الصلب head latency
بشكل عام ، كلما كان الـ head في القرص
اعلى ، في latency اكثر

to solve this problem, we have some tools
on windows and linux, to determine the
percentage of data scattering and then
doing some permutations to put the sectors
of the video or something directly next to
each other to improve the performance
(this is called fragmentation)

و ما هو الـ head في القرص الصلب و في القرص الصلب

~~* do not do fragmentation for SSD~~

ما هي الـ head في القرص الصلب ، كلما كان الـ head في القرص
اعلى ، في latency اكثر ، كلما كان الـ head في القرص
اقل ، في latency اقل

~~XXXXXXXXXXXXXXXXXXXX~~

Section 3 // cache memory

two questions! when reaching the cache to load:

- ① how to know if it really holds the data or not? (miss / hit)
- ② if it is (hit), how to find the data in the cache?

→ size of cache is much more smaller than DRAM

it is cache (between K-M bytes)
subset from DRAM (talking about Giga)

32 KB
64 KB
256 KB

II Direct Mapped cache:-

There are no lines here / blocks, ~~we~~ we just load the adjacent members (no usage of special locality principle) or temporal locality

cache:- (blocks structured)

each block capacity 1 word / 4 bytes

→ DRAM is also blocks-structured (more than cache)

Here in this design:

→ each location is a block, either in cache or

DRAM (4 bytes)

General rule to transfer the data from DRAM to cache:

II $\text{target address} \bmod \text{number of blocks in cache} = \text{the location to put the data in, in the cache}$

↳ This may cause conflicts

The conflict is described as: (many locations (words) in DRAM, they will map to the same location in cache → Direct Mapped cache)

[2] Use $\lceil \log_2(N) \rceil$

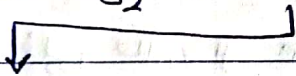
□ in cache, number of blocks is a multiple of 2

□ DRAM is the same

use $\lceil \log_2(N) \rceil$, N is the number of blocks in cache

so in slide 33, do the next:

① $\log_2(8) = 3$



② take the first three bits (least significant bits) from the right-hand-side, these bits will determine the location in cache.

This is done upon the assumption of:

each ~~4~~ ^{bytes} ~~words~~ are sharing the same block, and the same address in DRAM, so to load 4 word (4 bytes), you need one address.

is DRAM 4 bytes is still ~~same~~ DRAM address

→ This also can ~~also~~ cause a conflict

The question is

if one location in cache can hold the word of one of many locations in DRAM \rightarrow [conflict]
so how to determine the source?
of loaded data:

\rightarrow we used the ³ least significant bits, but the rest of bits is called the [Tag of Address] \rightarrow

so we load the data to cache, and the Tag to determine the source [cache] \rightarrow

\rightarrow We will also add a bit in cache $\left[\begin{matrix} 0 \\ \text{or} \\ 1 \end{matrix} \right]$ to determine if the block in cache is empty or not. (Valid bit)

→ it is necessary to modify the content of some block in cache, to be basically empty block, the index may be correct, but the tag is wrong, so we have the data of wrong location in DRAM, then we go to ~~and~~ request the right one from DRAM (because this is a miss), and modify the content of block + tag bit. This is because many blocks can be mapped to the same location in cache.

examples till slide 40: are for the cache that is word-based.

| | | | |
|----------------------|--------|--------|---|
| 4 bytes 32-bits ← | word 1 | 000000 | → each address is for a word |
| | word 2 | 000001 | |
| | word 3 | 000010 | → we use least significant 3 bits to identify the index |
| | word 4 | 000011 | the rest two are for the tag |

bytes ←
8 bits

| | | |
|--|----------|---------------------------------------|
| | 000 00 | → word 1 → first address of word 1 |
| | 000 01 | |
| | 000 10 | |
| | 000 11 | |
| | 0 0 1 00 | → word 2 → first address of word 2 |
| | 0 0 1 01 | |
| | 0 0 1 10 | |
| | 0 0 1 11 | |

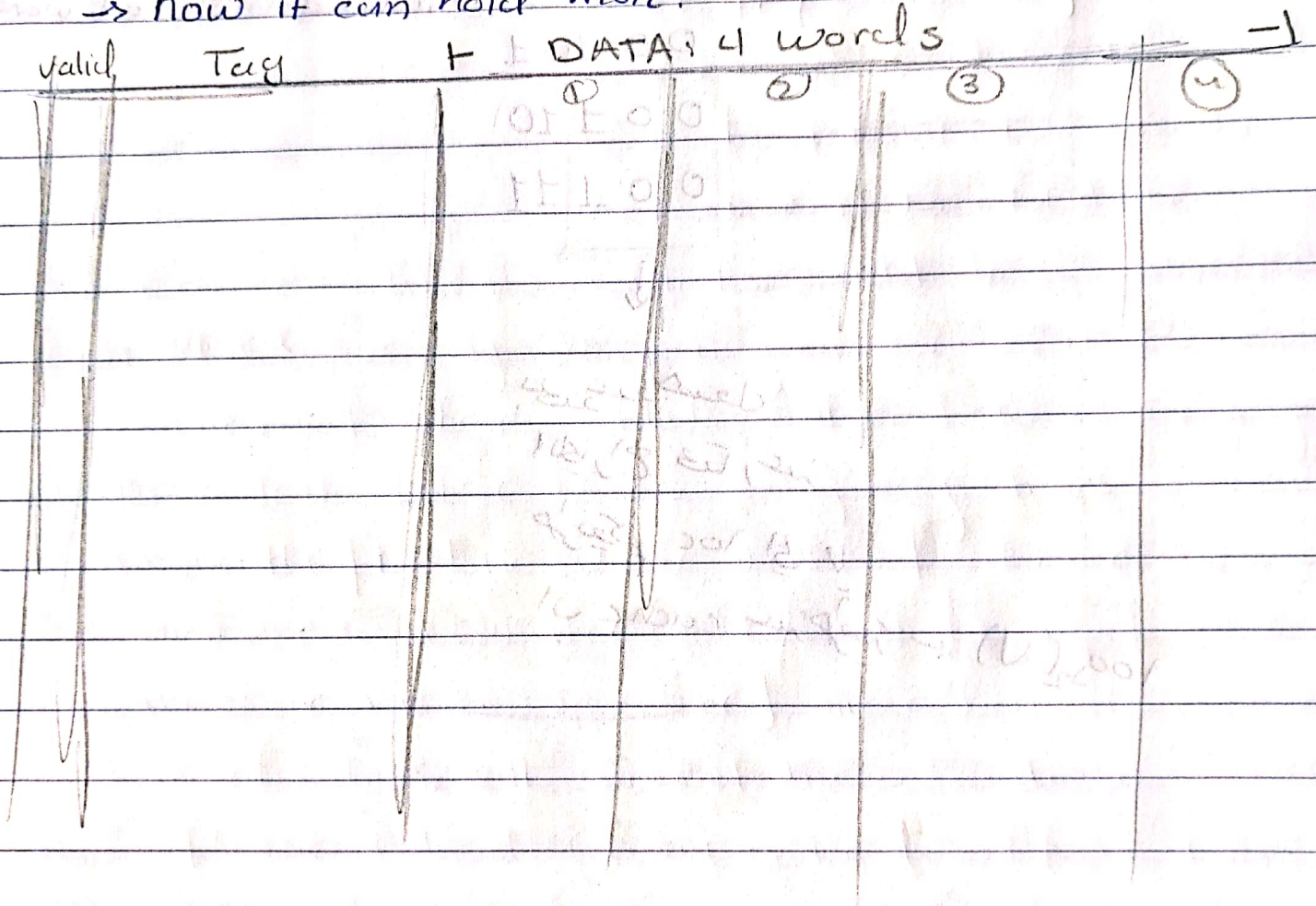


بفرض عدد هياكل
البرامح في الذاكرة
هو N و k هو حجم
البيانات في كل
cache $\rightarrow \log_2(N)$

Slide 42: A larger size cache

→ basically in one block it could hold one word

→ now it can hold more



هناك احتياجنا تقسيم indices حسب ال address ال موجودة
 مكونة من 4 words ، كل كلمة مع تقسيم ال offset ،
 bits

منه وبنه بيننا نكتبنا تقسيم ال bits الي يكون ال index

Slide 42:

The more ^{>big} line we get, the more words we load providing special locality for better performance but this true with some constraints only:

the constraint is the cache is small

after knowing that, know when the block of data becomes under than one word, we have two ways to identify where to save the data in the cache

1] Using mod in slide 42

2] Using the offset bits

First divide the address on the number of bytes in the block to determine the number of offset bits to know from where to start block to determine the group the block that you'll load and your target data will be in) →

then use mod

address mod
after
division

number of blocks
in cache

index of address of block

then load the block

Example: you have 32k bytes, design it in two ways: length of 256, length of 1024

II length of 256

the block size = $\frac{32k}{256} = 128$ byte [so wide] ↑
reduce miss rate

↳ This provides good spatial locality support, but 256 locations for a very big DRAM, will cause a lot of conflicts, a lot of addresses in DRAM are only competing on 256 locations in cache, and these conflicts will lead to modifying the content of these blocks many many times frequently, [128 byte will be modified at once when we load another address that compete on the same location in cache], these frequent modifications will conceal the advantage of spatial locality (why to load a lot of bytes while you know that you may modify them very much soon before you are able to access them in cache)

↳ Miss rate can increase here because of frequent changing of the content of blocks because of the competition on them

Block length of 1024

$$\text{size of block} = \frac{32 \times 1024}{1024} = 32 \text{ byte}$$

→ this solves the problem of the previous design (lot of conflicts, few locations) we have 1024 locations

→ if you are working on an array, you will get a miss each (8 elements / where the element is a word), not good handling for spatial locality

Miss rate can increase here because of small block sizes

→ this is trade off.

The question is on what to depend when

designing these caches?

→ design, simulate, determine what is the best for you!

What if we want to store to
data cache? [we can't write on
instruction cache]

Design in slide 46:

$$\text{number of offset bits} = \log_2(16) = 4$$

→ 10 bits for allocating

→ 18 (the rest) is for
tag

$$\text{management size} = 1024 \text{ rows} \left[\begin{array}{l} \text{one for} \\ \text{each} \\ \text{Block} \end{array} \right] * (1+1+18)$$

$$1024 * 20 = 20k$$

↓ bit

$$20 \text{ bits} \leftarrow 2.5 \text{ kByte}$$

- valid bit
- Dirty bit
- 18 bit of
tag

* What is the percentage of management
from the whole SRAM Block?

$$\frac{2.5 \text{ kB}}{(2.5 + 16) \text{ kbyte}} = 0.135$$

The D bit (Dirty bit)

→ we use it when we want to store (write)

The principle is, when you want to modify the value of some thing, modify it in the closest level to you (cache), not the RAM, because if we want to modify for more than one time, why to keep going to RAM in each time? just update the value on cache

↓
when you update the value of any bit in any location in cache, update its corresponding Dirty flag to 1, to announce that the value of this element is old in DRAM, and the newest is here in the cache

↓
The same variable has 2 values → old in DRAM
→ new in cache.

والبيته، cache الى الذاكرة Store

The most recent and all system also the up-to-date data is here in the cache subsystem

in Design slide 46#

→ when we loaded a block from DRAM to caches we loaded the target address + its adjacent

↓
but now how to transfer only the target data to CPU

جزء الذي يحتاجه البرمجيات ← check 26 بتات

بالنسبة إذا كانت من 32 بت إلى 64 بت
management bits

as a decoder bit enable

Question: on design 46 slide

if we revised the address from 32 bits to 64 bits, what is the size of management bits?

$50(\text{tag}) + 1(\text{D}) + 1(\text{V}) = 52 \text{ bits}$ [high overhead]

→ This is considered when talking about the cost.

Design of slide 46

Assuming we want to write on cache!

SD $\times 15$, $8(\times 31)$

Assumption (of current scenario)

- ① the address ^{in the} existing scenario is ~~the~~ is existing in cache to update its contents.

How cache works when it exists?

① calculated $2 \times 8 + (\times 31)$ to know the address

② identified the block that contains it and the word inside the block

The problem is you can't write ^{on only} (access only) one word [you either write ^{on one} and block or read one]

What to do after determining Block and word?

To store on cache you have two scenarios

- ① to load something from DRAM to cache
- ② to update some variable on cache from CPU coming data.

→ The overall scenario:

- ① data ~~is~~ read by [Data Read]
- ② all words read by it are taken to the muxes above.
- ③ the muxes above have 3 inputs:
 - A] input from old data read by [Data read]
 - B] input from the data coming from DRAM [Mem_data_data]
 - C] input from the CPU to update the data on cache [CPU_req_data]

Three cases: The address you are trying to access to store on may be ↓

Miss

↓
enable
Mem_data_data
input:

Hit

↓
enable
CPU_req_data
or
Data Read (old values) input

§ : means cache ✓

selection line on muxes identify which input to enable

→ we have a dedicated mux for each word.

When updating cache:

→ [Data Read] signals are sent on bus to muxes [each word to its corresponding mux]

→ CPU req. data is sent to each mux

→ The right mux is enabled only to upload the data to its word and update it.

→ rest of mux only pass the old values to words without any update.

→ The selection line determines which value to pass ~~through~~ to corresponding words.

~~Remember~~

→ Remember when writing to cache and updating it, we update the value of D bit to 1's (value is updated)

D is flipped to 1: DRAM is now inconsistent and synchronized with the values in cache to announce the update.

Question: what can we do instead of putting a dirty flag? to denote that there is an update.

↓
① do write-through update ~~in~~ cache and DRAM in the same time

Problem!! this way is slow

↓
trying to write on cache, 1 cycle or 2 are needed

↓
" " " " RAM, 100 or 200 cycles are needed.

II it is not acceptable to stop the program to update on cache and memory consuming all these cycles for a high delay [what is the advantage of cache then?]

② we may need to update the same value on cache multiple times after each other, it is not acceptable to go to RAM to update it each of these times

[time consuming]

↳ it is better to wait until you finish all these updates on cache then take the final value to RAM

The only advantage of write-through is the consistency [the system have has the most recent data every where]

*improvement on the design
(slide 48)

→ add the write buffer

→ the CPU can ~~still~~^{keep} working

you update the data in cache

and write it to ~~memory~~ write buffer

→ let the system load the updated

values from write buffer to Mem

shortly to be up-to-date

[CPU is still working, cache is not blocked]

Two problems:

[1] the more cost

[2] the small size of write buf.

assume the case of full filling the
write buffer, such that you can't write

on it, here we stall the cache [stop updating

on cache until we transfer data to

the buffer to RAM, and be ready again to

accept new updates]

The effective CPI, which means the penalty of store instructions, that works on a cache of write-through type, waits and writes on mem in the same time [the old write through] without the buffer, is 1

$$CPI = Old\ CPI + \frac{\text{percentage of store}}{\downarrow} * \text{cost of store (how much cycles it rates)}$$

(The new number of cycles needed to finish 1 inst. check slide 48)

→ With write buffer, we don't use the above equation except in the case that the buffer is full, such that we can't write on it any more.

write through is slow, buffer provides some improvement, but there is still some problems

Second way is to (Write Back)

→ note that both write-back and write through are being used, each is used for a ~~level~~ level. no one of them is better than the other.

dirty flag is not used in write through but used in write back

2) Write Back:

General Rule: do not update from cache to RAM until there is local miss (read miss) in cache so we need to write from mem to cache to handle that miss

↓

→ store instruction on cache to update location A

→ local instruction from cache to read from location

A

→ But location A has the data of different address than the target one

→ we want to load the right data from mem to location A

→ But we don't want to lose the updated values in location A

→ We store the updated value to mem, to ensure

we have the most recent data permanently, then we allow to load the missed data from mem to location A in cache to read it to CPU, while we are not worrying any more about our updated values

معالج ال local لا يقرأ من ال mem، و ال CPU

mem ال 100 cycle ←

cache ال 100 cycle ←

وال CPU ال 100 cycle

✓ To solve the problem of this delay, just change the sequence of execution:

instead of sending the update to mem, then loading data for local inst. from mem. to cache we do the next:

200 cycle

1) send the updated values from cache to temporary write buffer, and let the cache be ready to receive the data from mem to execute local

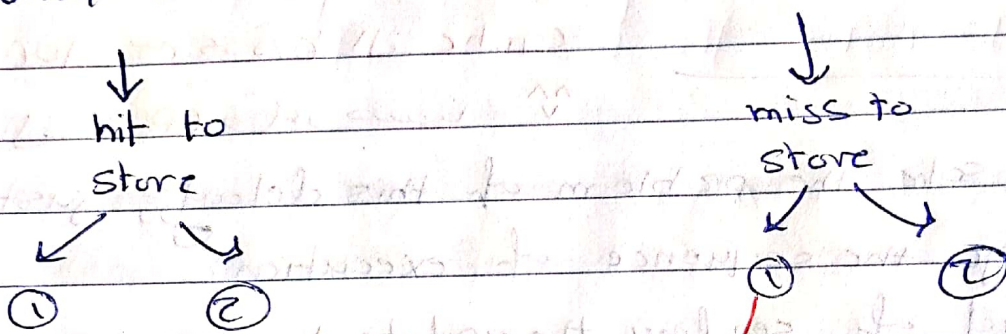
→ 100 cycles to transfer from mem. to cache, then sending them to the CPU, [CPU only waits 100 cycles instead of 200]

2) letting the write buffer to send the update to mem in any time shortly without forcing CPU to wait for this to be done ✓

The previous two scenarios indicate that for store operation (on cache), we always get (hit) - in other words, the address that we want to update its value is always existing in cache!

But what if this address is not existing!?
(miss)

what scenarios to be done // slide 51



① Write no Allocate:

you missed?, then your address is not in cache, just go directly to RAM and update it here.

you don't allocate, means you don't allocate a block and saves something on it

دفعه من الكاشه الى الذاكرة الرئيسية بدون تخصيص!

② write Allocate

→ local the data (cache missed address) from memory and its adjacent

→ put it in allocated block in cache

→ update their values of these local words

which is better? Allocate or no Allocate

→ no allocate is good for one variable, useless in spatial locality needs for arrays

→ allocate provides handling for spatial locality for things like arrays

it is true that we miss the first element in the array to update its value, but we are intending to update the values of the adjacent elements.

scenario: we want to update first 4 elements of an array!

→ no allocate:

miss $a[0]$ → 100 cycles to update in mem directly

" $a[1]$ " " " " " " "

" $a[2]$ " " " " " "

" $a[3]$ " " " " " "

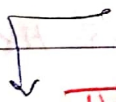
allocate:

miss $a[0]$? local $a[0]$ and its adjacent
and update them in cache with less cycles
required.



and only send these values to RAM
when some replacement is about to happen
because of [local miss] in the location that
holds the values of updated $a[0]$, $a[1]$, --

allocate ← loops للبرامج التي
no allocate ← التي



This question is answered by the statistical
results using ~~different~~ softwares, to determine
which way is better for each type of software,
if it is full of loops or something like that.

last five minutes in meta (23)

Weighted average in slide 54

→ The normal average is to take = $\frac{\text{I-cache miss rate} + \text{D-cache miss rate}}{2}$

This way assumes that each one is accessed 50% times of all accesses

- But in fact, this does not happen, we may have 20% I-cache accesses, for 80% D-cache accesses, so the weighted average is considering this: we multiply 0.2 * I-cache accesses and 0.80 * D-cache accesses

~~$$\frac{(0.20 * \text{I-cache miss rate}) + (0.80 * \text{D-cache miss rate})}{2}$$~~

Section 4# Measuring & improving cache performance.

Example to start: The program

[ideal case assumed where the data we want is always ready in cache (hit) & M stage takes one cycle.]

LD

LD

ADD

SUB

SD

MUL

MUL

LD

DIV

SD

SUBI

ADDI

ADD

SD

LD

MUL

SD

1) determine the percentage of inst. that access the mem. to the number of all inst. [these are LD + SD]

$$\text{Memory Accesses} = \frac{8}{17} \approx 0.47$$

Program

CPI \downarrow scheduling cycles \downarrow inst. \downarrow 47%

2) determine if all of these LD and SD are having penalty of miss or not

وال IPC \downarrow But how

what is the number of inst. that have a miss?!

the cache misses affect

miss \leftarrow inst. \downarrow miss rate

[the base or ideal CPI

miss rate

of ideal cases (always

3) Multiply the miss penalty to miss rate to find the overall miss penalty

hit)]

In English: from the entire program we find the instructions that are accessing memory $\left[\frac{\text{mem accesses}}{\text{program}} \right]$, then we find the inst. that have misses from these instructions [miss rate], then we find the overall misspenalty to find **memory store cycles** // slide 57

real example // slide 58

in these equations of this example, they indicate the point of view of the pipeline, for know about its ability and how many cycles it can has to take to finish an inst.

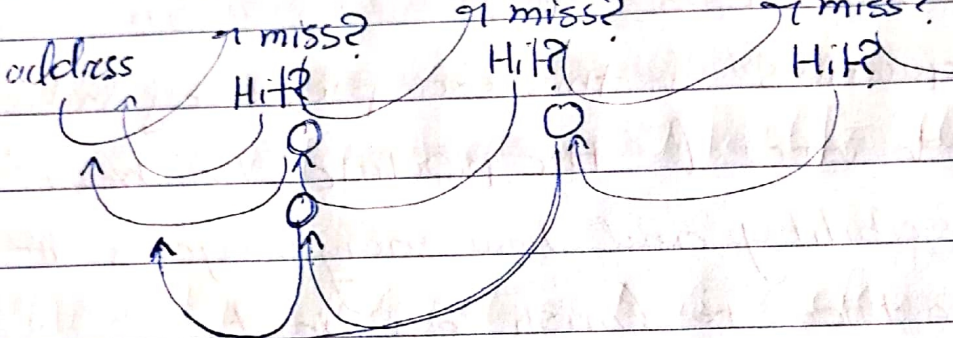
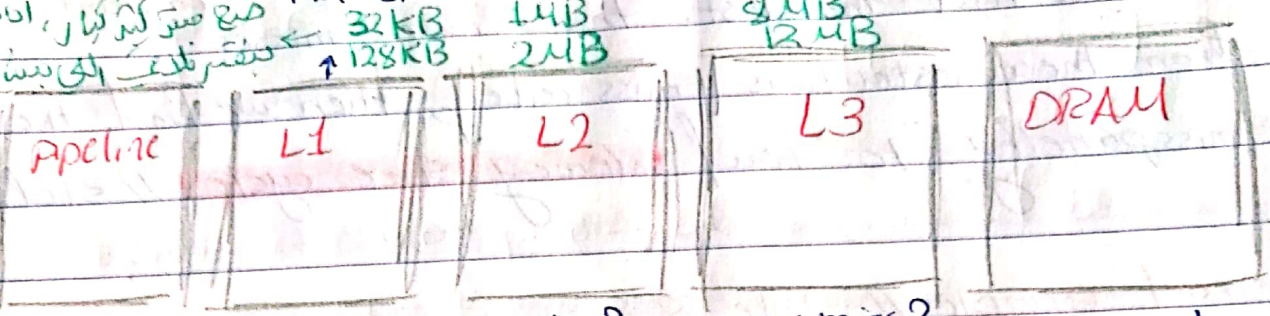
> slide 59: use the point of view of the time [how much time we need to load our data from RAM to handle the miss]

More realistic examples on that with more levels

of caches in stack 60,

their sizes are not high

مع صغر حجم الكاش، انحصار دونه
تقتصر تلك التي تبينها
لا يوجد
لا يوجد
كثرة



الاصول التي تبنى
فيها، وخصيص
زوج عال Hard
disk
تكونت نسبة ال

- Hit in level 1: → pipeline
- " " " 2: → L1 → pipe
- " " " 3: → L2 → L1 → pipe

ال disk لها
تت سطون سن
على ال RAM

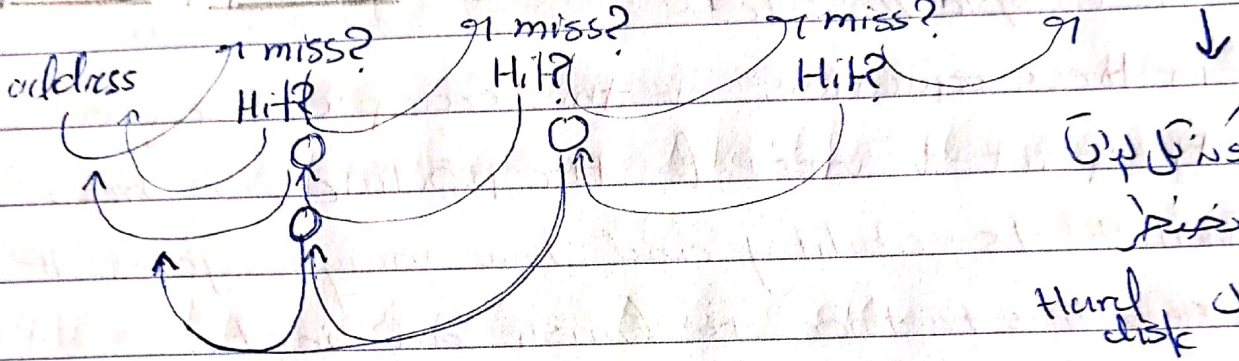
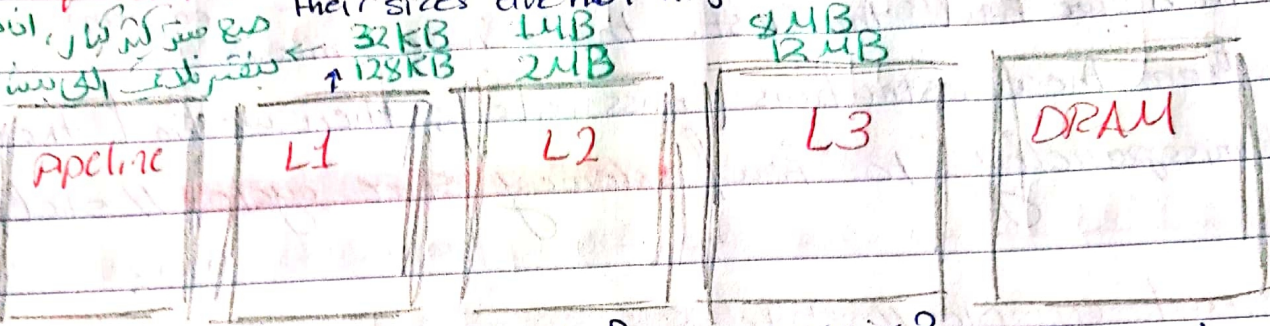
and so on

→ we pay the complete penalty when we should reach the DRAM level because of all misses

More realistic examples on that with more levels

of caches in state 60, their sizes are not high

مع صغر حجم الكاش، انشغال الكاش بغير تلك التي نريدها
 لا يصح
 بعض البيانات
 كبيرة



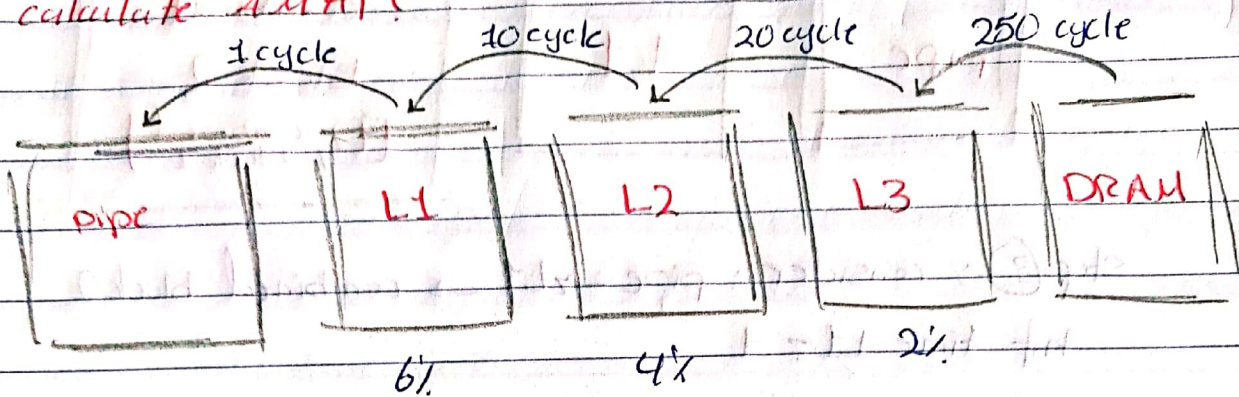
الأصل يكون في الذاكرة
 فيها ما نريده
 نرجع على
 Hard disk
 من ذاكرة الذاكرة

- Hit in level 1: → pipeline
- " " " 2: → L1 → pipe
- " " " 3: → L2 → L1 → pipe
- ;
- and so on

ال disk لها
 في سطر
 على ال RAM

→ we pay the complete penalty when we should reach the DRAM level because of all misses

In slide 60: how to analyze the system to calculate AMAT?



AMAT = hit time + miss rate * penalty

eg unit: $1 \text{ cycle} \times \text{pipe} + 6\% \times 10 \text{ cycles} + 4\% \times 20 \text{ cycles} + 2\% \times 250 \text{ cycles} + \text{DRAM}$

how to apply it to a hierarchical system?

step 1: consider: pipe \rightarrow L3 \rightarrow DRAM

ignore 2 levels: pipe, L1 level

hit time L3 = 20 cycles

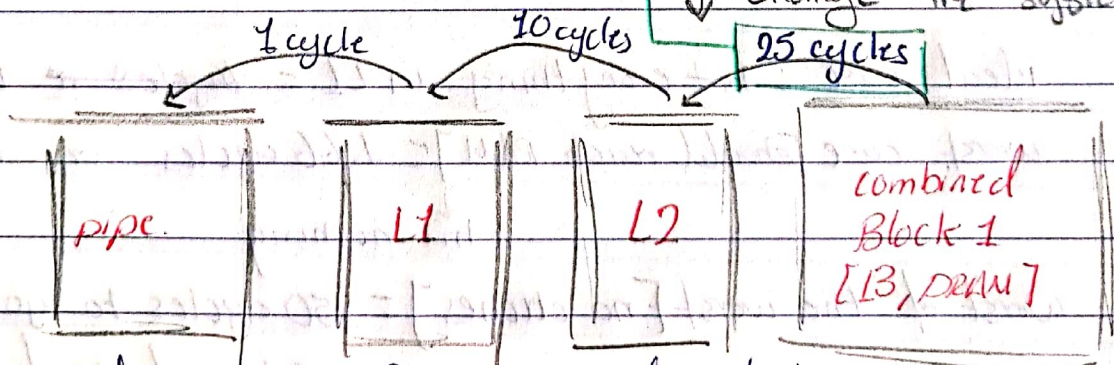
miss rate L3 = 2%

penalty to load from DRAM to L3 (Miss penalty) = 250 cycles

operation always start from the slower level towards the pipeline

$$AMAT = 20 + 0.02 * 250 = 25 \text{ cycles} \rightarrow AMAT[L3, DRAM]$$

change the system



step 2: consider: pipe \rightarrow L2 \rightarrow combined block 1

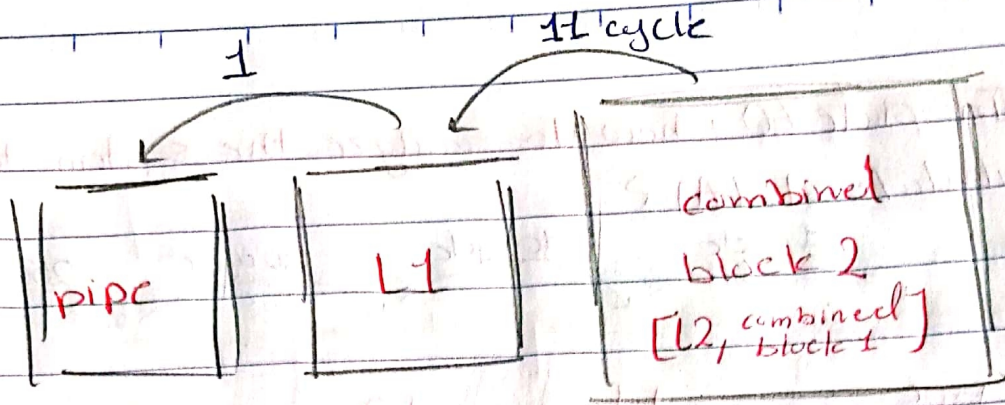
hit time L2 = 10

miss rate L2 = 4%

$$\rightarrow AMAT = 10 + 0.04 * 25 = 11 \text{ cycles}$$

penalty (combined to L2) = 25

change the system



step ③: consider: pipe \rightarrow L1 \rightarrow combined block 2

hit time $L1 = 1$

miss rate = 6%

penalty [combined block 2 to L1] = 11 cycles

AMAT = $1 + 0.06 \times 11 = 1.66$ cycles

So Average Memory Access time from CPU view for cycles and times $\checkmark = 1.66$ cycles

[Analyze this] number: 1.66 \checkmark

ideal case: hit everything in L1 = ~~cycles~~ 1 cycle

worst case [should reach RAM] = 1.66 cycles *on average*

hit nothing

worst of the worst [no caches] = 250 cycles to go to

RAM and local data

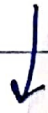
from it

ideal \rightarrow 1.66, 250 \rightarrow cache \rightarrow P₂, P₃

in the measurement.

→ the more we added levels of caches the more we could reach ~~more~~ better AELAT, and closer to the ideal

This lead to the point of: Build large and fast memory hierarchy, but the two traditional designs we could have are:



Fast but small

[pipe ↔ cache]



large but slow

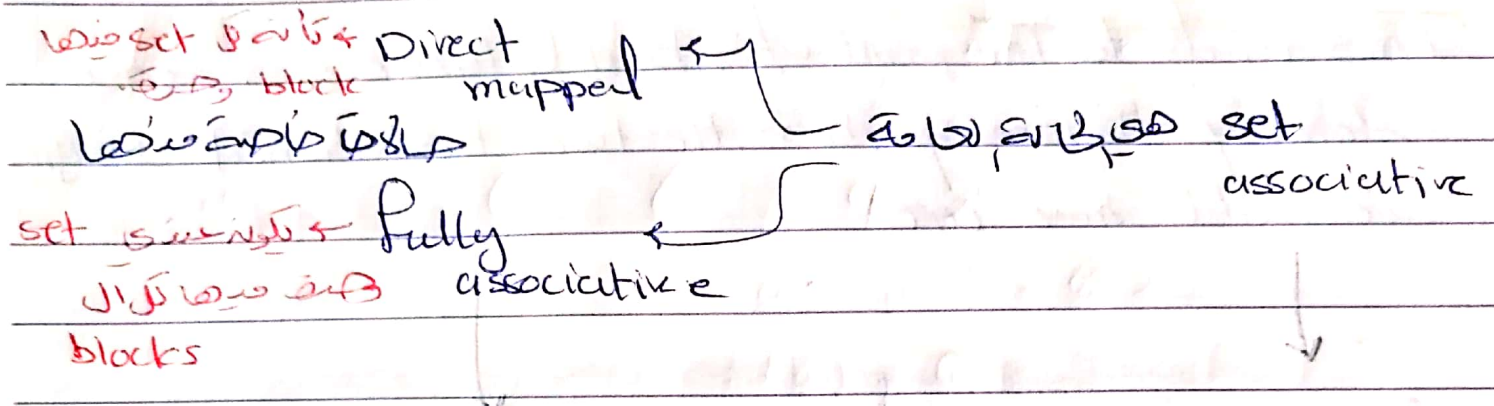
[pipe ↔ DRAM]

But, the hierarchical design of caches [levels of caches] could lead us to the point of large memory, but ^{with} fast performance, with some specifications like the example

slide 63

* Set-associative

instead of dividing cache to blocks, divide to sets, and each set contains ^{a group of} ~~some~~ blocks

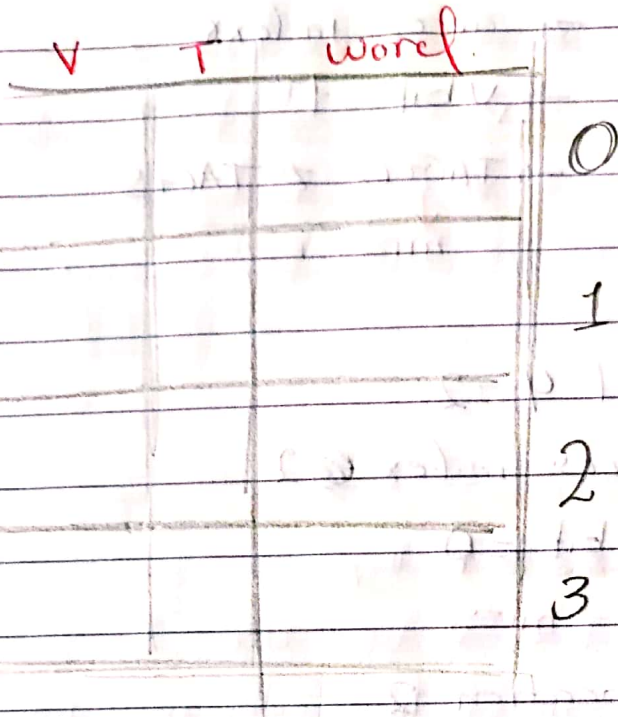


Example:-

- 2-way set associative cache
- 8 ~~blocks~~ only locations
- How many sets? 4 sets

what if it is 4-way? 2 sets per 8 ~~blocks~~ locations

slide 65: the cache:



because only 4 blocks we have.

① first address: $0 \bmod 4 = 0$

cycle #1

→ access cache index 0

→ cache is empty (V flag = 0)

go to RAM ← ↓ miss ✓

and load wanted data to location 0

② 2nd address: $8 \bmod 4 = 0$

→ access cache index 0

→ V bit = 1

→ Tag of 8 \neq Tag of 0

go to RAM ← ↓ miss ✓

and load data to location 0

③ address 0 again: $0 \bmod 4 = 0$
→ access index 0
→ V bit = 1
→ TAG 0 \neq TAG 8
↓
miss ✓

④ address 6 : $6 \bmod 4 = 2$
→ access index 2
→ V bit = 0

local data from ← miss
RAM and put it to location 2

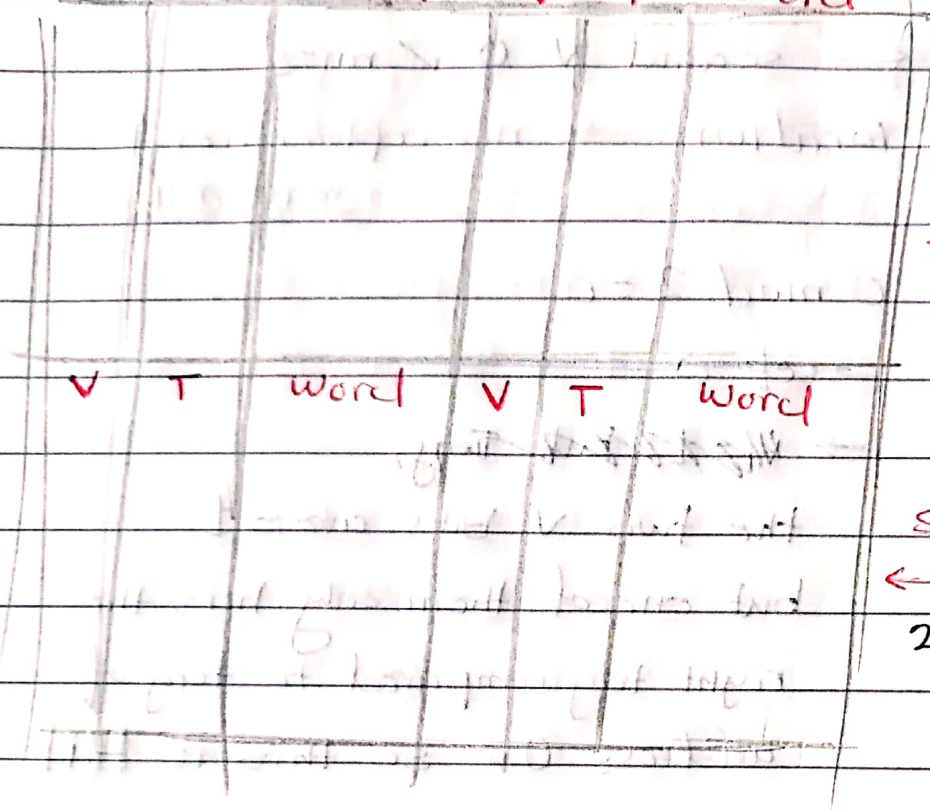
⑤ address 8 : $8 \bmod 4 = 0$
→ index 0
→ V bit = 1
→ TAG 8 \neq TAG 0
miss ✓

hit rate for this example for
Direct mapped? hit rate = 0%

slide 68 :

- 2-way set associative

V T word V T word



Set 0

← 2 words

Set 1

← 2 words

because we have 2 sets: we will divide on the number of sets in the equation, not the number of blocks [equation of mod]

① address 0: $0 \bmod 2 = 0$

→ access set 0

→ compare the two V bits in the set

↳ both are = 0 (cache empty)

go to RAM, and

← miss ↓

load the wanted data, to one of these words in the set (any one of them)

② address 8 : $8 \bmod 2 = 0$

→ set 0

→ compare first $V = 1$ ✓

miss ← but $\text{Tag } 0 \neq \text{Tag } 8$

load from RAM ← second $V = 0$ ✓ miss

to this empty location → no replacement required

③ address 0 : $0 \bmod 2 = 0$

→ set 0

→ ~~the two V bits~~ ~~are = 1~~

the two V bits are = 1

but one of them only has the

right tag compared to tag of

address 0 ✓ so this is **Hit**.

④ address 6 : $6 \bmod 2 = 0$

→ set 0

→ $V = 1$

$V = 1$

load from RAM

← $\text{Tag } 6 \neq \text{Tag } 0$

$\text{Tag } 6 \neq \text{Tag } 0$

and we will assume

miss ✓

to replace the least recently used

data (0, 10) → This is the location of Mem[8]

replace it!

⑤ address 8 : $8 \bmod 2 = 0$

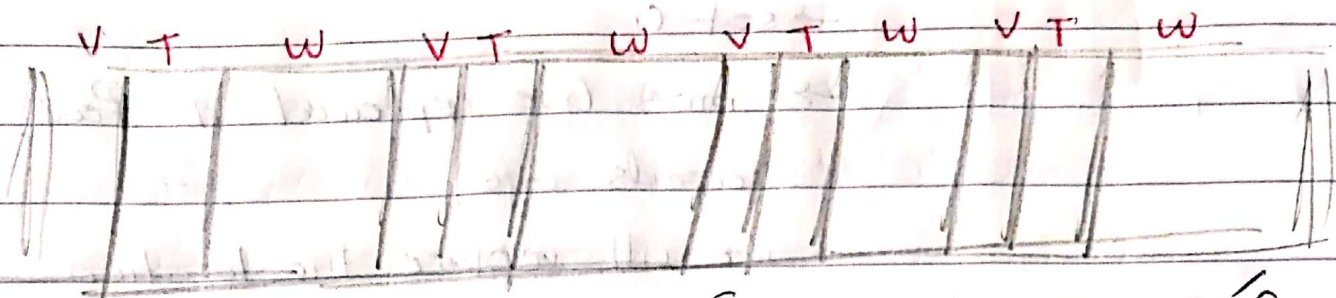
→ set 0

→ miss (we replaced it few seconds ago)

→ we will replace the location containing address 0 (the least recently used)

$$\text{hit rate} = \frac{1}{5} = 0.2 = 20\%$$

slide 66: Full associative



تسلكنا من طرف او قطبة و طرف للبا: V, T, W بجيب V, T, W في V

no need to use mod // no indexing or sets

(1) address 0: \rightarrow compare all V bits in cache (all are = 0)

miss \checkmark

\rightarrow load data from RAM to Randomly the first location \checkmark

(2) address 8: \rightarrow first $V=1$, Tag 0 \neq Tag 8

\rightarrow rest of V bits are = 0

\rightarrow load from RAM to randomly 2nd location \checkmark

(3) address 0: \rightarrow first $V=1$, Tag is correct \checkmark

Hit \checkmark

④ address 6: \rightarrow 1st $V=1$, Tag 0 \neq Tag 6
 \rightarrow 2nd $V=1$, Tag 8 \neq Tag 6
miss $\leftarrow \rightarrow$ rest of V bits are $=0$
 \rightarrow load from RAM to the first empty place (3rd location)

⑤ address 8: \rightarrow 1st $V=1$, Tag 0 \neq Tag 8
 \rightarrow 2nd $V=1$; Tag is correct \checkmark
hit \checkmark

$$\text{hit rate} = 40\% = \frac{2}{5} \times 100\%$$

Law of diminishing returns

مثلاً ، لو بدنا نضيف على عمال يستغلوا حرفة لينجروا
تسجل محسن ، اذا اقبلنا نزيد عددهم لمحسن ، بصير تسجل
اسرع بدون تاثير سلبي ، بس اذا زدنا عددهم عندهم كحد
بصير تاثير لزيادة سلبي وعكته تكونه ~~تكونه~~ تكثري

او بصير تكبير السرعة الاجاز تسجل حتى كثير الة تاثير
واضع (تكونه حافيف) ، وهاد مثال الكلفة لمجموعة
على عدد العمال

design slide 68:

→ checking the bits of V and Tag of all locations
in one set is done in parallel

→ because from each set we only load one
word (one location), we don't need the block offset
(we don't load the entire line)

→ 256 sets : needs 8 bits to identify the index of
set

→ 22 bits are left for Tag

The Byte offset is related to the structure of DRAM, and how it is built.

→ when the structure was horizontally, (4 bytes next to each other, we didn't need to use the byte offset)

→ and when the DRAM is 32: word is 32 bits (4 bytes) byte offset = 2

" " " " = 64 : " " = 3

↳ word is 64 bits (8 bytes)

In the design of Direct Mapped:

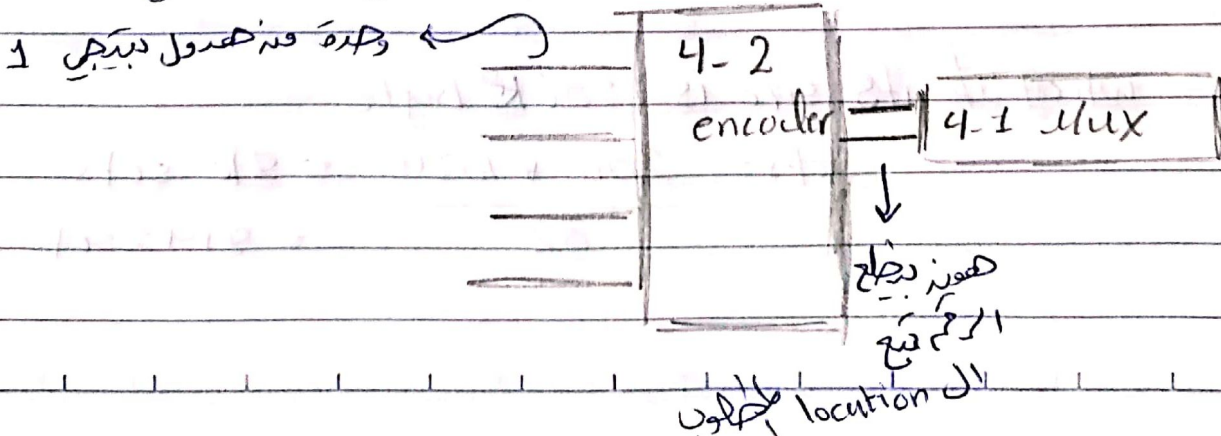
we used the block offset to supply it to the Mux that chooses the line (as a selector line)

→ But here in set associative we don't have a block offset, ~~not~~ ^{how} to identify the location needed?

→ We need 2 lines to select

But in design of slide 68: 4 lines are getting into the mux as the selection line, why?

it may be wrong or, we have an encoder



Assume the design is for fully associative, will find differences and cost!

- there will only be one line // as we have one set
- mux will be wider to handle all these locations at the same time

Example:-

cache → 8-way // each location is of 4 word
1 word = 4 byte
set has 4 byte × 8 = 32 bytes

cache size 32k byte
how many sets

$$\text{sets numbers} = \frac{32 * 1024 \text{ byte}}{32 \text{ byte}}$$

we have 1024 sets.

what if its size is 256k byte

$$\text{sets} = \frac{256 * 1024}{32} = 8192 \text{ sets}$$

what if its size 6MB and it is 12-way?

width is : $12 \times 4 = 48$ byte for each set

$$\text{sets } \frac{6 \times 1024^2}{48}$$

why L3 is shared between all cores? slide 72

Two reasons!!

[1] Avoid Contention, what would happen if each one of them had its own L3 cache?

→ Basically the CPU is connected to the bus, to the DRAM, and because L3 cache is the last level then it's the one that CPU requests will pass to reach DRAM, then L3 is connected to the bus, then if each core had its own L3, then each one will have a connection point on the bus (shared bus), so contention will happen, because if a lot of cores tried to request something through the bus, L3 caches will compete to win the connection and pass the signal to DRAM

② parallel processing, such that a lot of cores will work on one array for example in parallel, so it is better to keep it in a shared level like L3 cache, to keep up-to-date ~~with~~ with the updates coming from ~~for~~ other cores on this array.

Interactions with software!! slide 77.

elements and arrays are stored in two different ways in memory:

- row major (bytes & elements are stored in row)
- column major (bytes & elements are stored in columns)

row like C++ and Java

column, MATLAB

→ Misses depend on memory access patterns (How the code affect the performance?)

Example: on C (row major)

```
for (int i=0; i<N; i++) {  
    for (int j=0; j<M; j++) {  
        arr[i][j] = rand() % 10; }  
}
```

outer loop iterates over rows
inner loop iterates over columns

access `arr[0][0]`: miss in cache

load `arr[0][0]`: and a lot of its adjacents due to special locality, then a lot of elements will be hit

This code benefits from this special locality because it is row major (iterates over rows) → the code is written to do that. and the language is row major itself.

~~How?~~

- How? - hardware is row major // identical to language
- when cache loads data from DRAM, it loads the adjacents in rows
 - language is row major
 - code is written to iterate over rows
 - cache has the adjacent elements of the row
 - very low miss rate

If one modification happened on code to iterate over columns not rows: `arr[j][i] = rand() % 100;`

we won't benefit the special locality, because in cache we have adjacents of a row, not a column, all elements will be missed!

To test the cache effectively and the DRAM, and test how much hits and misses are affected and reacting with code??

we have to use kind of large array such that it does not fit completely in level 3 cache so we are forced to go in some cases to DRAM + make the number of iterations big to some level such that it satisfies the same idea

The experiment of testing the system, is repeated for thousands of times to get the average of all results, which is better.

Section 11 // Redundant Arrays of Inexpensive Disks

Harddiskse ۱۱ // Redundant Arrays of Inexpensive Disks

The trend now is going to use cloud instead of RAID systems:

- for less electricity usage
- " " " cooling systems need
- " " " cost of paying for RAID physical devices

* But, it is not the best choice for private data (security)

→ parallelism increases bandwidth, for which you need a part of your data from each disk, which increases performance [this is done in parallel instead of taking all chunks of data serially one after the other, which takes more time]

Hot swapping?

have made electrical designs of interfaces that allow the hot swapping, such that you can remove one hard disk from the system, and replace it with another, without needing to shut down the whole system for a long time, which makes it down and not responding for frequent accesses [This is a problem]

Types of raids due to different designs:

RAID 0, RAID 1, RAID 2, ...

Main difference between RAID 2 and RAID 3

[2]

- split at bit level
- parity is stored on multiple disks

[3]

- split at byte level
- parity is stored on one disk only

Main difference between RAID 3, and RAID 4

B)

- splitted at byte level

F)

- at block level

blocks have big sizes
which is satisfying

RAID 4 // Native

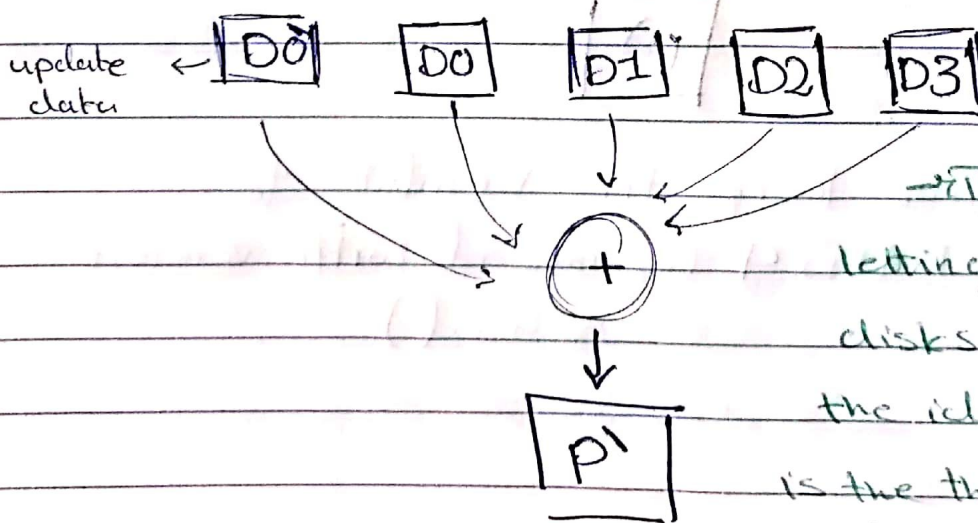
→ When to update data:

□ D0' is put instead of D0

But parity disk should be updated too ✓

↓

Reuse all disks to recalculate P' through XOR circuit with the new data



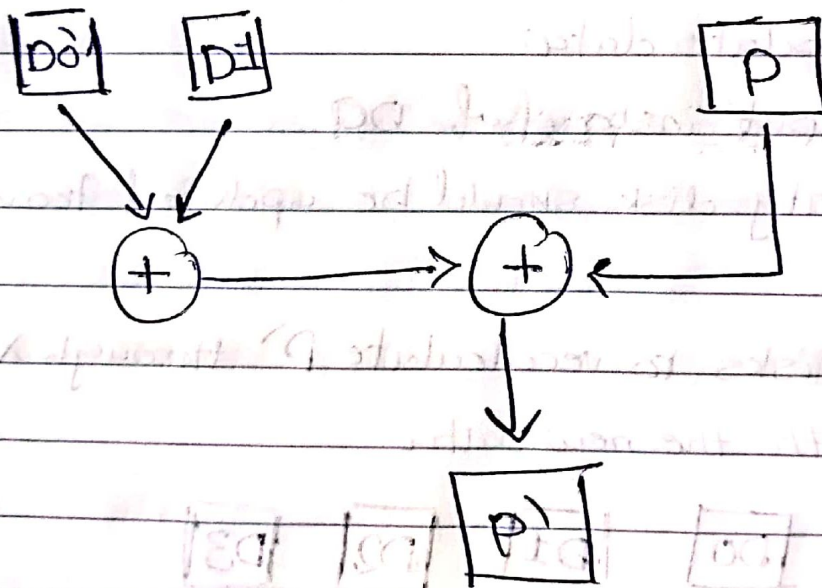
→ This requires letting all the disks work and leave the ideal mode, which is the thing that we tried to avoid by designing RAID 4 !!

This have a high cost on disks!!
all benefits are gone!!!

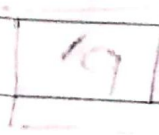
RAID 4 // Optimize write

→ no need to use the ~~completed~~ ^{un-updated} disks to recalculate the parity

you will have the same result by applying XOR between old data and new data, then apply XOR between the result of first XOR and old parity



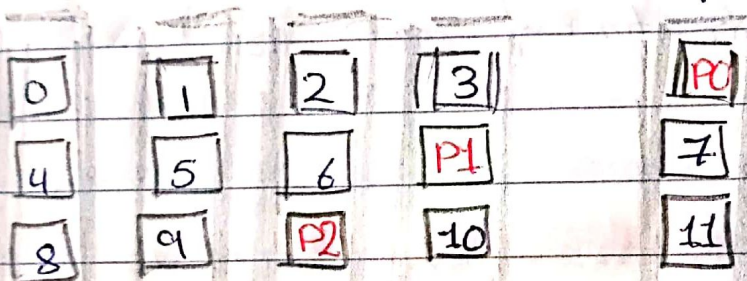
This way, we keep the benefit of designing RAID 4 (optimized write version is used)



upon ~~all~~ all the benefits of RAID 4 design, ~~we~~ we have 2 new benefits come with RAID 5 design.

□ if the parity disk failed, and got down for some reason

the ~~data~~ diagonally distributed ~~chunks~~ chunks of parity will find the lost data in parity disk



P0 is lost
data 7 is lost
" " "
- you still can
retrieve data
7 by
P1, 6, 5, 4
- and retrieve
data 11 by
P2, 10, 9, 8 ✓

→ knowing that in some designs like RAID 3, if any one of data disk has failure, we can retrieve all its data using the rest of data disks and parity disk

But the problem here, that if the failure disk is the parity disk, the case when you can't retrieve any thing!

[2] we should not always use the parity disk to update data or retrieve it

~~data~~

parity units are distributed to many disks, so whenever you need to use ECC for any reason you are not forced to access the parity disk necessarily which [if you are forced to do so] could decrease the life time of parity disk!

→ Problem of ~~RAID~~ RAID is

we store all the system in one place, ~~because~~

so when some physical damage happen to

that place, the entire system will be affected

to solve this problem, use Multiple RAID ✓

section 6# Virtual Machines!!

advantages of virtual machines:

[1] testing alot of programs on different operating systems, only on one physical device (development advantages)

[2] Good for cloud computing, where you can use a shared hardware that provides the cloud and runs different VM for different OS of users [security + reliability]

↓
using a VM
can isolate
you safely
and
completely
from other
users or
OSs

↳ when the OS
stops responding
suddenly, you
can solve the
problem by only
shutting down the
VM and restarting
it again
while other users
& operating systems
and VMs are still
working

These days VMs are supported by the hardware itself to accelerate it.

↳ Virtual Machine Monitor

* The VMM does the job of mapping translation from the logical tasks of guest OS and real physical execution that should be done.

* User mode: Apps

* privileged mode: instructions dedicated for OS

The problem is the Guest OS running on the device in user mode, while it is an OS and should be using the instructions of privileged mode.

when it tries to do so, exception is thrown, and this exception is raised/trapped to VMM as a request to be allowed to use the protected resources.

↑

This is another reason to dedicate parts of HW for virtual machines and guest operating system.

* if one core is logically dedicated for the
Guest OS, then the real physical execution
can be entirely different after translation
by VMM

Section 7# Virtual Memory

ۋاڤا قىلىۋېتىش

* Cache is managed by the physical gates of hardware

virtual machine
ۋاڤا قىلىۋېتىش

But when working with RAM, OS should be involved in the work to manage it with hardware support

Contents

- 5.1 Introduction
- 5.2 Memory Technologies
- 5.3 The Basics of Caching
- 5.4 Measuring and Improving Cache Performance
- 5.5 Distributed Memory Hierarchy
- 5.11 Redundant Arrays of Inexpensive Disks
- 5.6 Virtual Machines
- 5.7 Virtual Memory
- 5.8 A Common Framework for Memory Hierarchy
- 5.9 Using a Finite-State Machine to Control a Simple Cache
- 5.10 Cache Coherence
- 5.13 The ARM Cortex-A53 and Intel Core i7 Memory Hierarchies
- 5.16 Fallacies and Pitfalls
- 5.17 Concluding Remarks



RAID

Redundant Array of Inexpensive (Independent) Disks

- ← Use multiple smaller disks ^{instead of} (c.f. one large disk)
 - Parallelism improves performance ✓ *[Because some data has copies distributed to all these disks]*
 - Plus extra disk(s) for redundant data storage
- Provides fault tolerant storage system →
- Especially if failed disks can be "hot swapped"

if one of them fails, you need to change only that unit but not that whole entire system

لو فشل واحد بقدر
سببه عطل

↓
The modern
standard.



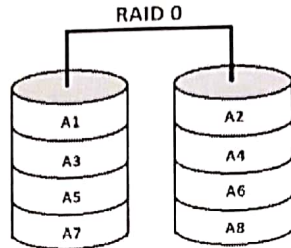
RAID 0

→ Basically there is no RAID system } only a combination of ~~hard disk~~ hard disk

No redundancy ("AID"?)

Just stripe data over multiple disks → But it does improve performance

internally: the data is distributed in way such that each chunk of data is stored in a different hard disk (There are no copies)



advantages:
1) reading data in parallel improves performance.

disadvantages:
→ chunks of data can be lost with no Backup



RAID 1

advantage
1) dependability; Backup for data
2) less performance because you are forced to read serially always from one hard disk (data disk)

RAID 1: Mirroring

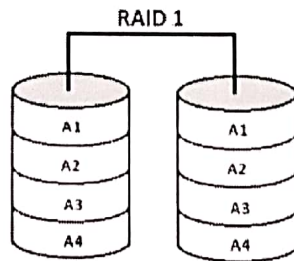
N + N disks, replicate data

Write data to both data disk and mirror disk

* On disk failure, read from mirror
determined after RAID detection

* you only have access to data disk for writing or reading

* Gives you a time to replace the failing disk through hot swapping and after replacing it internally, all the data updated on the mirror disk will be copied to the new data disk



Have n number of hard disks anything written on one hard disk will be internally copied to other hard disks.

so for example:
if a system is containing 2 of 6TB hard disks, then the system size will be 6TB, because you store 6TB on one of hard disks and the rest of them are Backup hard disks you can't save more than 6TB



starting from the idea that fault tolerance should include things more than dependability (Backups) and parallel access, things like flipping bits inside disk (which harms the data) should be solved

↳ [Mechanically it is working, but the data is not correct] → logical tolerancy
 ↳ store Error correction code with Data [ECC]

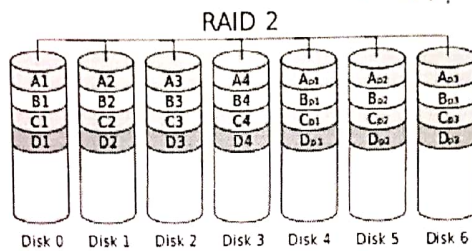
RAID 2

so RAID 2 has become

RAID 2: Error correcting code (ECC)

▪ N + E disks (e.g., 10 + 4) N for data disks, E for parity disks

- Split data at bit level across N disks
- Generate E-bit ECC
- Too complex, not used in practice



Too complex because:

- it is complex to split at bit level
- bit level data requires you to let all the disks work to provide data.
- (high cost at all levels)
- no one of them can stay idle for a small period of time

Chapter 6 — Storage and Other I/O Topics — 96

This increases the rate of disks fail to correct the wrong data.

- when to read:
- read all disks together
 - regenerate ECC
 - compare it with stored ECC
 - if some thing went wrong



it corrects it

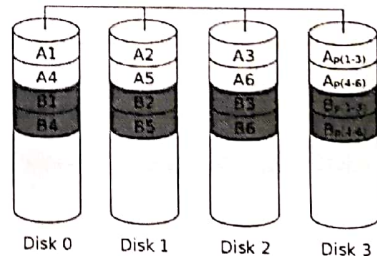
- the provides you with the data requested.

RAID 3: Bit-Interleaved Parity

N + 1 disks

- Data striped across N disks at byte level
- Redundant disk stores parity
- Read access
Read all disks
- Write access
Generate new parity and update all disks
- On failure
Use parity to reconstruct missing data

instead of bit level



صلى الله عليه وسلم
 RAID2 نية
 نية

Not widely used (still a little bit complex)



Chapter 6 — Storage and Other I/O Topics — 97

→ Have weak points because of:

- 1) bytes splitting is still complex (not much better than bit level)
- 2) storing parity in one disk does not provide any dependability advantage so that parity can be lost at failure

correct the wrongy
cluster.

RAID 3: Bit-Interleaved Parity

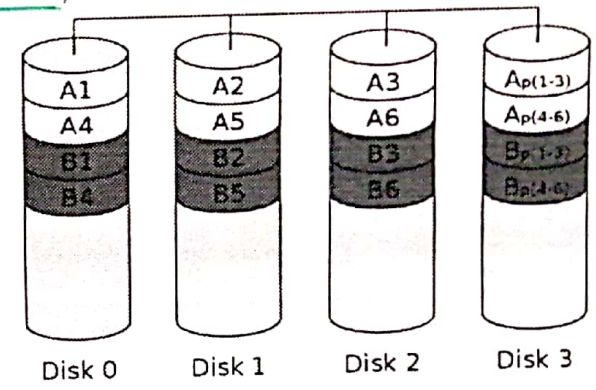
N + 1 disks ^{parity disk}

- Data striped across N disks at **byte level**
- Redundant disk stores parity

instead of bit level
RAID 3

صحيح كذا
RAID2
نوعه

- Read access
Read **all** disks
- Write access
Generate new parity and update **all** disks
- On failure
Use parity to reconstruct missing data



Not widely used (still a little bit complex)



→ Have weak points because of:

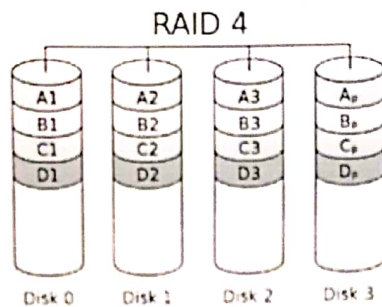
- 1) bytes splitting is still complex (not much better than bit level)
- 2) storing parity in one disk does not provide any dependability advantage so that parity can be lost at failure and you lose the protection

RAID 4: Block-Interleaved Parity

N + 1 disks] like RAID 3

- Data striped across N disks at **block level**
- Redundant disk stores parity for a group of blocks
- Read access

Read **only** the disk holding the required block



parity is calculated at block level.

mechanically better for hard disk to increase the time it will remain able to write

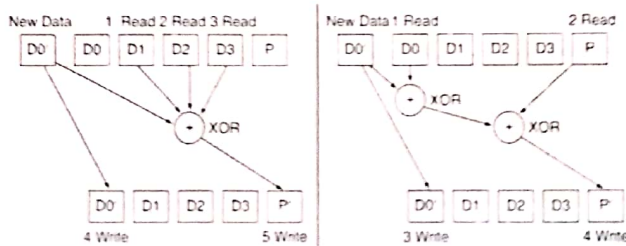
Big sizes of blocks that are distributed to disks, allows that when you request some data, it all can be of size less than that block or equal to it, so you are not forced to let all the disks work to provide you the data, each time you are accessing some thing

two designs !!

RAID 4: Naive vs. Optimized Write

- Write access
 - Just read disk containing modified block, and parity disk
 - Calculate new parity, **update data disk and parity disk**
- On failure
 - Use parity to reconstruct missing data

Not widely used



| A | B | XOR |
|---|---|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

disadvantage:

that comes about ~~fault~~ ^{Disk corruption} fault tolerance is gone
 having the parity ~~entirely~~ entirely stored in one disk, what makes it hard to retrieve it when disk fails (data will be lost), so we lose the protection on data

RAID 5: Distributed Parity

improvement on RAID 4 to solve the problem of storing the parity on one disk

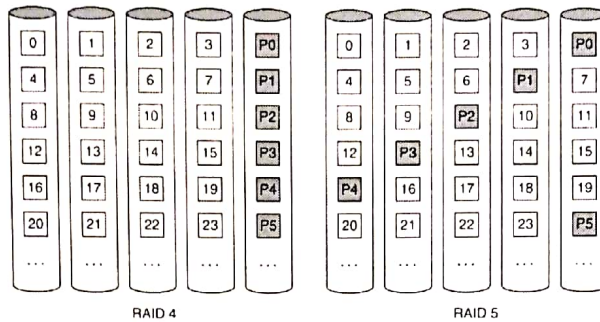
N + 1 disks

- Like RAID 4, but parity blocks distributed across disks
- Avoids parity disk being a bottleneck

Widely used ✓

Mixing

distributed in a diagonal way!



RAID 4

RAID 5



RAID 6: P + Q Redundancy

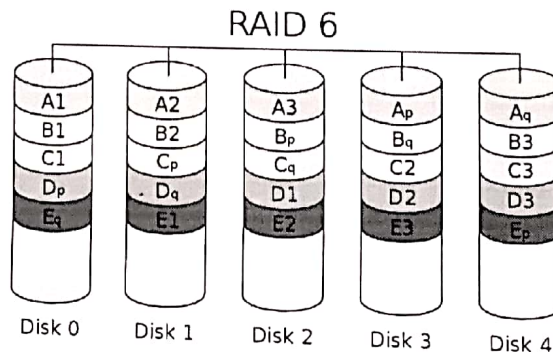
Same as RAID 5 but with using more than one block for parity, which is also distributed diagonally.

N + 2 disks

- Like RAID 5, but two lots of parity
- Greater fault tolerance through more redundancy



stronger protection because of redundancy



Disk 0

Disk 1

Disk 2

Disk 3

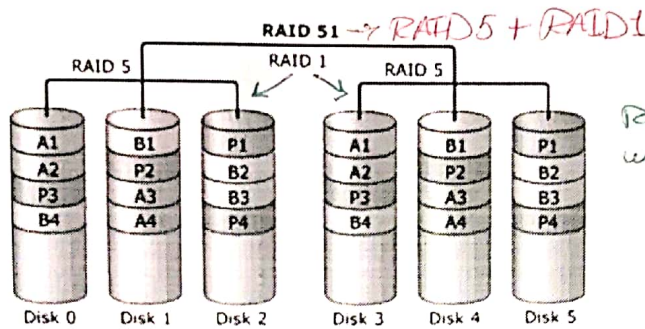
Disk 4



Multiple RAID

Multiple RAID

- More advanced systems give similar fault tolerance with better performance
- Example RAID 51



RAID 51 is to do replication with no parity generation

all the system here is copied to another system in different geographically place

Locally: RAID 5
Geographically: RAID 1

RAID Summary

RAID can improve performance and availability

- High availability requires hot swapping

Assumes independent disk failures

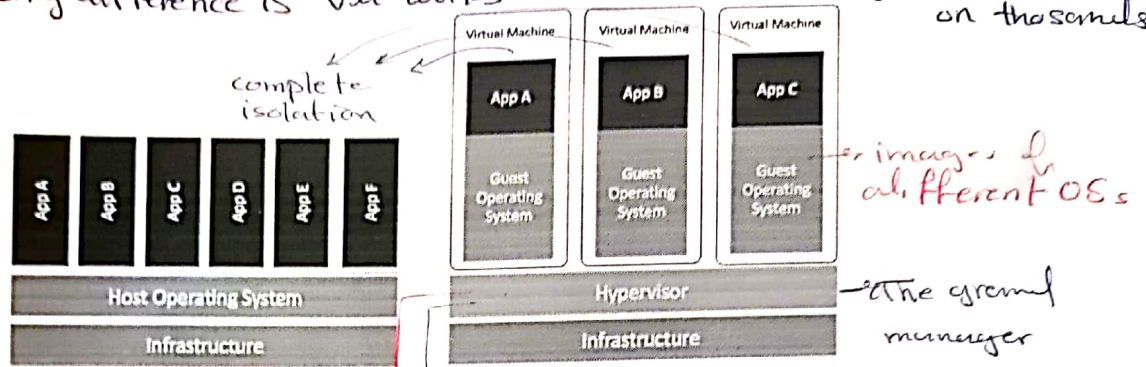
- Too bad if the building burns down!

summary in lecture 28

| | Data disks | Redundant check disks |
|---|------------|-----------------------|
| RAID 0 (No redundancy) Widely used | 4 disks | 0 disks |
| RAID 1 (Mirroring) EMC, HP(Tandem), IBM | 4 disks | 4 disks |
| RAID 2 (Error detection and correction code) Unused | 4 disks | 3 disks |
| RAID 3 (Bit-interleaved parity) Storage concepts | 4 disks | 1 disk |
| RAID 4 (Block-interleaving parity) Network appliance | 4 disks | 1 disk |
| RAID 5 (Distributed block-interleaved parity) Widely used | 4 disks | 1 disk |
| RAID 6 (P + Q redundancy) Recently popular | 4 disks | 2 disks |

Virtual Machines

VM and Hypervisor are almost the same
the only difference is VM works on the device, while Hypervisor works on thousands



run each App for a user on a Guest OS, inside a VM

The only problem is only when the server of Hypervisor gets down

* The single operating system (the Guest) does not have complete access on the host machine specifications completely, it just can access what you allowed it to access

Virtual Machine Monitor

Maps virtual resources to physical resources

- Memory, I/O devices, CPUs

→ don't dedicating resources for that Guest operating system

Guest code runs on native machine in user mode

- Traps to VMM on privileged instructions and access to protected resources

Guest OS may be different from host OS

VMM handles real I/O devices

- Emulates generic virtual I/O devices for guest