# Content

# Eight Great Ideas

ملاحظة بين إنه كل سنة أو سنة ونص
بنضاعف عدد الـ transistors اللي بنقدر
نحطهم على 1 integrated circuit

مؤسس شركة intel
- Design for *Moore's Law*

الأشي المعقد تخليه أبسط (تخلي الأسهل هو المكشوف ف ليسهل التعامل).
- Use *abstraction* to simplify design

الحالة اللي بنستخدمها كثير نخليها أسرع. (ما نسرع الكل عشان التكلفة والقلة)
- Make the *common case fast*

تقسيم الشغل على أكثر من شخص (توزيع الشغل على أكثر من طريق)
- Performance *via parallelism*

الشغل على أكثر من مهمة على التوالي
- Performance *via pipelining*

الشغل بناء على الشيء المستقبلي
- Performance *via prediction*

لما أبني الـ memory ما بخليها شيء واحد (نقسمها)
- *Hierarchy* of memories

قريش الحسابات اللي بعملها بنثق فيها و فيه بفمد و بوثق فيه
- *Dependability via redundancy*

بطونا
high
perf..
بنقدر ننجز
أكثر من اشي

ما بنحط كل شي بالـ cash
1: كل ما أكبر الشي بصير أبطئ
2: التكلفة عالية جداً

لو وجود بديل لو خرب اشي في الـ
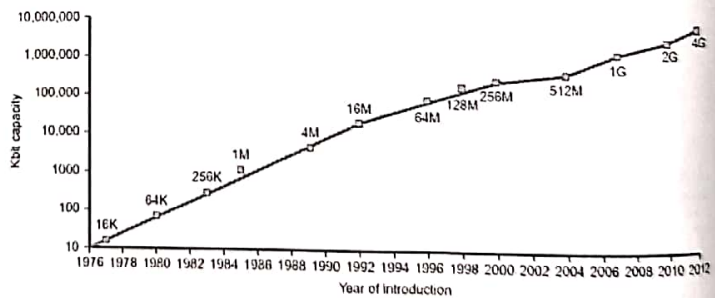بديل مثل وجود مروحتين تبريد

Scanned with CamScanner

# Content

بحتاج لتحقيق higher Performance: وجود تكلفة و دوائر لعمل المهمات
و بعالم الحواسيب هاي الدوائر لازم تكون أكبر و بتكون أغلى و أسرع.

# Technology Trends

\* transistor فيها أكثر من chip اللي بينما الواحد transistor في هي Vaccum tube \*

- Electronics technology continues to evolve
  - Increased capacity and performance
  - Reduced cost



DRAM capacity

| Year | Technology | Relative performance/cost |
|------|-----------|---------------------------|
| 1951 | Vacuum tube | 1 |
| 1965 | Transistor | 35 |
| 1975 | Integrated circuit (IC) | 900 |
| 1995 | Very large scale IC (VLSI) | 2,400,000 |
| 2013 | Ultra large scale IC | 250,000,000,000 |

Scanned with CamScanner

# Semiconductor Technology
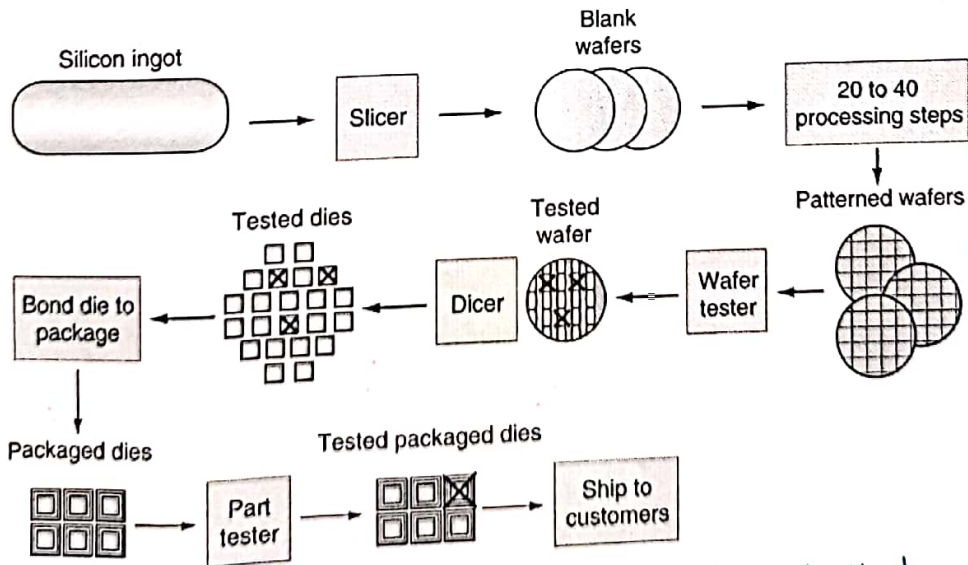
ما المادة الرئيسية وفوقها

Silicon: semiconductor

Add materials to transform properties:

- Conductors موصل
- Insulators عازل
- Switch مفتاح

# Manufacturing ICs

" كيف بتشتغل "

" مكلفة عشان هيك حاليا شركات معينة اللي بتعملوا بس "
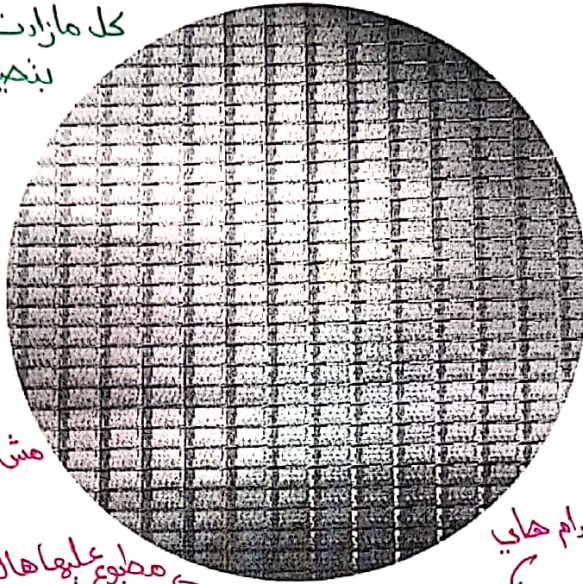


* كل ما بكون عندي خبرة عدد القطع الخرابة بتقل فالـ Yield بصير أعلى فالتكلفة بتقل و ما إلى ذلك .

← عدد الـ dies الصالحات تقسيم العدد الكلي

- Yield: proportion of working dies per wafer نسبة

لأنه بالـ wafer ممكن يكون فيه dies معطلة .

# Intel Core i7 Wafer

كل مازادت مساحة الـ die
بنصير تكلفتها أعلى

مش كلام ممكن يكونوا احالجين ← مطبوع عليها القطر

بنحدد قديه الـ transistor صغار أوكبار.

باستخدام هاي

- 300mm wafer, 280 chips, 32nm technology
- Each chip is 20.7 x 10.5 mm

أبعادها

كلما كان الرقم أصغر
الـ transistor بكونوا أقل.

( cost ∝ area ← ولكن مش، ممكن يكون تربيعي linear ) يعتمد

# Integrated Circuit Cost

الـ cost العلاقة بين

مساحة مستطيل م ← مساحة الدائرة

$$\text{Cost per die} = \frac{\text{Cost per wafer}}{\text{Dies per wafer} \times \text{Yield}}$$

عدد الـ chips ← Dies per wafer ≈ Wafer area/Die area
بالـ wafer

بتطلع 280 تقريبا
chips

$$\text{Yield} = \frac{1}{(1 + (\text{Defects per area} \times \text{Die area}/2))^2}$$

منطقة فيها خربان transistor

★ كلما كانت الكثافة أقل كلما كان الـ Yield أفضل

لـ بالأخط أكبر من 280
بس لأنه اللي
جايين عأطراف
الدائرة مش كاملين
فبنلغيهم

ثابت
Nonlinear relation to area and defect rate

- Wafer cost and area are fixed

تحدد      معدل الخلل
- Defect rate determined by manufacturing process

- Die area determined by architecture and circuit design

★ Defect Per area = $\dfrac{\text{# of defects}}{\text{wafer area}}$

# Content

# Response Time and Throughput

- Response time
  - How long it takes to do a task → قديه وقت لأنفذ برنامج أو instruction
- Throughput
  - Total work done per unit time → قديه بنقدر نعمل شغل بوحدة الزمن
    - e.g., tasks/transactions/... per hour
- How are response time and throughput affected by
  استجابة
  - Replacing the processor with a faster version?
  - Adding more processors?
- We'll focus on response time for now...

# Relative Performance

- Define Performance = 1/Execution Time
- "X is *n* time faster than Y" « $Performance = \frac{1}{T}$ »

لح لهنا يكون أكبر ما يمكن

ب

لهنا يكون

أقل ما يمكن

عدد ال Parallel units

$Performance_X / Performance_Y$

ي مشواحد لا P ↑

$Throughput = \frac{P}{T}$

= $Execution\ time_Y / Execution\ time_X = n$

Time العملية ←

الواحدة

لـ لما اشتغل عال Parallelism أو Piplining مش دقيقة فبهيس →

Example: time taken to run a program

- 10s on A, 15s on B
- $Execution\ Time_B / Execution\ Time_A$
  = 15s / 10s = 1.5
- So A is 1.5 times faster than B

# Measuring Execution Time

- Elapsed time
  - Total response time, including all aspects
    - Processing, I/O, OS overhead, idle time النوائي
  - Determines system performance
- CPU time

الوقت المستغرق في معالجة وظيفة معينة

- Time spent processing a given job
  - Discounts I/O time, other jobs' shares
- Comprises user CPU time and system CPU time
- Different programs are affected differently by CPU and system performance

# CPU Clocking

- Operation of digital hardware governed by a constant-rate clock

يتحكم



Clock period

Clock (cycles)

Data transfer and computation

Update state

- Clock period: duration of a clock cycle
  - e.g., 250ps = 0.25ns = $250×10^{-12}$s
- Clock frequency (rate): cycles per second
  - e.g., 4.0GHz = 4000MHz = $4.0×10^9$Hz

# CPU Time

$$CPU\ Time = CPU\ Clock\ Cycles × Clock\ Cycle\ Time$$

$$= \frac{CPU\ Clock\ Cycles}{Clock\ Rate}$$

- Performance improved by
  - Reducing number of clock cycles
  - Increasing clock rate
  - Hardware designer must often trade off clock rate against cycle count

Scanned with CamScanner

# Instruction Count and CPI

Clock Cycles = Instruction Count × Cycles per Instruction

قديه فيه instruction          و زمن ال cycle الوحدة

CPU Time = Instruction Count × CPI × Clock Cycle Time

↳ clock per instruction

$$= \frac{\text{Instruction Count} \times \text{CPI}}{\text{Clock Rate } (f)}$$

$f = \frac{1}{T}$ ← Clock Rate $(f)$

- Instruction Count for a program
  - Determined by program, ISA and compiler
- Average cycles per instruction
  - Determined by CPU hardware
  - If different instructions have different CPI
    - Average CPI affected by instruction mix

وقت لتخلص task

Ex: Time = 0.1 Seconds / CPI = 0.5 / IC = 400 million Instructions

لا يقدر ينفذ 2 instru في ال cycle الوحدة          ( f = ?? )

Sol: $0.1 = 400 \times 10^6 \times 0.5 \times T$

$\to T = 5 \times 10^{-10}$

$f = \frac{1}{T} \to f = 2 \times 10^9$

$= 2$ GH!!

# CPI in More Detail

- If different instruction classes take different numbers of cycles

$$\text{Clock Cycles} = \sum_{i=1}^{n}(\text{CPI}_i \times \text{Instruction Count}_i)$$

- Weighted average CPI

$$\text{CPI} = \frac{\text{Clock Cycles}}{\text{Instruction Count}} = \sum_{i=1}^{n}\left(\text{CPI}_i \times \frac{\text{Instruction Count}_i}{\text{Instruction Count}}\right)$$

Relative frequency

أهم معادلة

# Performance Summary

*كل ما خل ال time زاد ال performanceال*

IC عدد instruc.. في البرنامج

$$\text{CPU Time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock cycle}}$$

الوحدة اللي بتطلع → CPU time

= قديه محتاج Seconds لينفذ Program معين

- Performance depends on
  - Algorithm: affects IC, possibly CPI
  - Programming language: affects IC, CPI
  - Compiler: affects IC, CPI
  - Instruction set architecture: affects IC, CPI, $T_c$

MK

Chapter 1 — Computer Abstractions and Technology — 18

كيفية تقليل:
1 - $\frac{\text{Instructions}}{\text{Program}}$ : بيعتمد عالخوارزميات و كيف استخدموا للتقليل من الـ Instructions بالعملية الوحدة.

2 - $\frac{\text{Clock cycles}}{\text{Instructions}}$ : نخلص أكثر من Instruc.. بنفس ال cycle يعني نعمل Pipelining وهكذا!

# Content

## 1.2 Eight Great Ideas in Computer Architecture
(Review) high frequency يعني علشان لازم اشتغل على اقل ال Period لو بدي اقلل ال

3 - $\frac{\text{Seconds}}{\text{clock cycle}}$

عن طريق استخدام Smaller transistors لانهم بصيروا اسرع

## 1.5 Technologies for Building Processors and Memory

## 1.6 Performance (Review)

## 1.7 The Power Wall

## 1.8 The Sea Change: The Switch from Uniprocessors to Multiprocessors
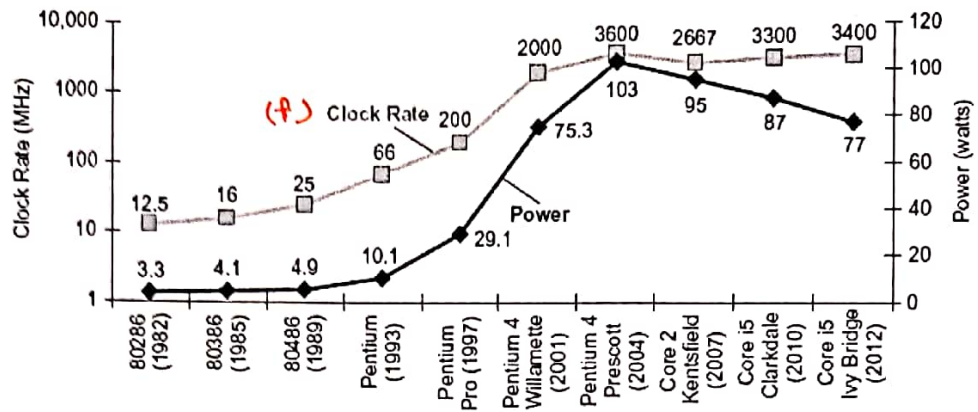
## 1.9 Real Stuff: Benchmarking the Intel Core i7

## 1.10 Fallacies and Pitfalls

## 1.11 Concluding Remarks

# Power Trends

حنحكي عن استهلاك الطاقة الكهربائية

\* الـ Processors بستهلكوا طاقة كهربائية وهاي الطاقة مع الزمن بتزداد \*



## In CMOS IC technology

الـ transistors الي فيه ما في تيار بمر فيهم فبننظر الام إنهم capacitors الي بيشبعوا

بشحن و يفرغ

الـ sur نتناسب طرديا معاها

Dynamic Power

$$\text{Power} = \frac{1}{2} \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency}$$

( يعتمد على : حجم العدد الـ transistors )
كل ما قل ر بقل الـ Power حجم الـ transistor

×30  يعني كل ما زاد الـ transistor عددهم

5V → 1V

×1000  لـ بتزداد الفترة

لـ من 5 ل 1 وهلا أقل من 1 شوي وما بنزل أكثر من هيك
له بتزداد الفترة

بقل حجمهم بقل حجمهم بقل الـ Power

\* الـ static Power الـ chip دايما الـ بتستهلكها \* = supply voltage \* current leakage Power

\* بطلنا نزيد الـ Frequency عشان نقلل الـ time لأنه كل ما زادنا الـ Frequency بتزداد الـ Power
فصرنا ندور حلول أخرى لتقليل الـ time

# Reducing Power

- ## Suppose a new CPU has
  - ■ 85% of capacitive load of old CPU
  - ■ 15% voltage and 15% frequency reduction

$$\frac{P_{new}}{P_{old}} = \frac{C_{old} \times 0.85 \times (V_{old} \times 0.85)^2 \times F_{old} \times 0.85}{C_{old} \times V_{old}^2 \times F_{old}} = 0.85^4 = 0.52$$

- ## The power wall
  - ■ We can't reduce voltage further  ما بنقدر ننزله أكثر من هيك
  - ■ We can't remove more heat  ما بنقدر نشيل remove لأكثر من 100 أو 150 watt
  - ## How else can we improve performance?

لـ بزيدوا الـ performance بالوقت الحاضر عن طريق الـ Multiprocessor

# Content

MK MK

بالا chip الوحدة

\# of logical cores لأنه بطلنا نقدر نقلل الا voltage أقل من هيك ، بس الحل اجون بزيادة الـ

قدرتنا على تحسين الـ performance لـ one processor نتقل بالوقت الحاضر وصلت 3.5%

# Uniprocessor Performance



Constrained by power, instruction-level parallelism, memory latency

بزداد بـ rate أقل   22%/year

الـ cpu الوحدة كانت تتحسن بـ rate عالي جدًا

52%/year

25%/year

# Multiprocessors

- Multicore microprocessors
  - More than one processor per chip
    لازم نلجأ ال Multicore لنستفيد من
- Requires explicitly parallel programming
  - Compare with instruction level parallelism
    حل في عبئ على المبرومجين انهم لبعدوا كتابة برامج Parallel program حتى تستفيد من ال Multiple cores
    - Hardware executes multiple instructions at once
    - Hidden from the programmer
- Hard to do
  لانه برامجك لازم تكون مح ولازم تقطيلها Performance أعلى
  - Programming for performance
  - Load balancing نخليه قلبل وسريع
  - Optimizing communication and synchronization
    ده لانهم بحتاجوا وقت ليتفاهموا مع بعض

بدك تقسم العمل على أجزاء متساوية بحيث الدورة ع Multiple
كل core ياخد شغل قد التاني لأنه لو واحد أخد شغل أكثر من غيره رح ياخذ كل البرنامج ولازم بيشوه لحتى يخلص

# Content

# SPEC CPU Benchmark

"طرق قياسية تستخدم للمقارنة بين جهازين أو تصميمين وهكذا."

- Programs used to measure performance
  يمثلوا ال Workload الحقيقيات ثم ال Benchmark
  - Supposedly typical of actual workload →
- Standard Performance Evaluation Corp (SPEC)
  مختلفة
  - Develops benchmarks for CPU, I/O, Web, …

- SPEC CPU2006 → فعالة وفي قبلها كان وفي بعدها 2017
  - Elapsed time to execute a selection of programs
    SPEC بتركز ع CPU ← Negligible I/O, so focuses on <u>CPU</u> performance
  - Normalize relative to reference machine → $\frac{T_{Ri}}{T_i}$
  - Summarize as geometric <u>mean</u> of performance ratios
    CINT2006 (integer) and CFP2006 (floating-point)
    هو الي بيستخدمه → Summerising Perf.

Selection of program
يصلوا الى 43
لم اختيارهم لأرقم

$$\sqrt[n]{\prod_{i=1}^{n} \text{Execution time ratio}_i} \rightarrow \frac{T_R}{T_i} \rightarrow$$ older reference machine
وقت ال → $T_R$
الوقت تاعك → $T_i$

Typical of actual workload مضبوط

م يعني → CPU integer

# CINT2006 for Intel Core i7 920

كلهم نشتاه لجسبوا $T_i$ →

| Description | Name | Instruction Count x 10^9 | CPI | Clock cycle time (seconds x 10^-9) | Execution Time (seconds) | Reference Time (seconds) | SPECratio |
|---|---|---|---|---|---|---|---|
| Interpreted string processing | perl | 2252 | 0.60 | 0.376 | 508 | 9770 | 19.2 |
| Block-sorting compression | bzip2 | 2390 | 0.70 | 0.376 | 629 | 9650 | 15.4 |
| GNU C compiler | gcc | 794 | 1.20 | 0.376 | 358 | 8050 | 22.5 |
| Combinatorial optimization | mcf | 221 | 2.66 | 0.376 | 221 | 9120 | 41.2 |
| Go game (AI) | go | 1274 | 1.10 | 0.376 | 527 | 10490 | 19.9 |
| Search gene sequence | hmmer | 2616 | 0.60 | 0.376 | 590 | 9330 | 15.8 |
| Chess game (AI) | sjeng | 1948 | 0.80 | 0.376 | 586 | 12100 | 20.7 |
| Quantum computer simulation | libquantum | 650 | 0.44 | 0.376 | 109 | 20720 | 190.0 |
| Video compression | h264avc | 3793 | 0.50 | 0.376 | 713 | 22130 | 31.0 |
| Discrete event simulation library | omnetpp | 367 | 2.10 | 0.376 | 290 | 6250 | 21.5 |
| Games/path finding | astar | 1250 | 1.00 | 0.376 | 470 | 7020 | 14.9 |
| XML parsing | xalancbmk | 1045 | 0.70 | 0.376 | 275 | 6900 | 25.1 |
| Geometric mean | - | - | - | - | - | - | 25.7 |

Integer app
وبرمجيات مختلفة.
نفس كردم م
بضربهم ببعض كلهم بعدين بياخذ الجذر ال 12 لالهم
كل الناتج

# SPEC Power Benchmark

Power consumption of server at different workload levels

- Performance: ssj_ops/sec
- Power: Watts (Joules/sec)

Power benchmark

$$\text{Overall ssj\_ops per Watt} = \left(\sum_{i=0}^{10} \text{ssj\_ops}_i\right) \Big/ \left(\sum_{i=0}^{10} \text{power}_i\right)$$

مجموع ↗   مجموع ↗

* The SPEC cpu 2017 benchmark package contains 43 benchmarks, organised into four suites:

1- SPEC speed integer and SPEC speed floating ] are used for comparing time for a computer to complete single tasks

2- SPEC rate integer and SPEC rate floating ] measure the throughput or work per unit of time

بعمني لو بدي أشوف الجهاز
قديه بياخذ وقت للمهمة
الوحدة

بعمني لها بعي أشتغل أكثر من
مهمة بنفس الوقت

نستخدم بـ general app مثل compiler أو Sorting أو Exel
وهكذا. بستخدم للبرمجيات اللي بتعملنا أمور مختلفة

غالبية الحسابات العلمية والهندسية تستخدمها.

# SPECpower_ssj2008 for Xeon X5650

operation per second ← وحدته

| Target Load % | Performance (ssj_ops) | Average Power (Watts) |
|---|---|---|
| 100% | 865,618 | 258 |
| 90% | 786,688 | 242 |
| 80% | 698,051 | 224 |
| 70% | 607,826 | 204 |
| 60% | 521,391 | 185 |
| 50% | 436,757 | 170 |
| 40% | 345,919 | 157 |
| 30% | 262,071 | 146 |
| 20% | 176,061 | 135 |
| 10% | 86,784 | 121 |
| 0% | 0 | 80 |
| Overall Sum | 4,787,166 | 1,922 |
| Σssj_ops/Σpower = | | 2,490 |

بشتغل Processor الـ على different target load

لـ كل ما كان أكبر كلما كان أفضل

كل ما زاد الـ target load بنجد الـ average power بزيد في بشتغل وكل مازاد : لأنه بحير في بشتغل وكل مازاد بزيد الاستهلاك للطاقة.

# Content

قديم Computer eng ملخص

# Pitfall: Amdahl's Law

- Improving an aspect of a computer and expecting a proportional improvement in overall performance

نقسم ال exec. time
الى Taffected
و Tunaffected

الزمن الجديد

$$T_{improved} = \frac{T_{affected}}{improvemen\ t\ factor} + T_{unaffected}$$

بنعدل عليه

ما بنقدر نعدل عليه

Instructions

- Example: multiply accounts for 80s/100s → بياخدوا وقت 80 من أصل 100

  - How much improvement in multiply performance to get 5× overall?

$$20 = \frac{80}{n} + 20$$

  - <u>Can't be done!</u>

$5 = \frac{T_{old}}{T_{improved}}$  speedup

$5 = \frac{100}{T_i}$

→ $T_i = 20$

- Corollary: make the common case fast

Scanned with CamScanner

# Fallacy: Low Power at Idle

- Look back at i7 power benchmark
  - At 100% load: 258W
  - At 50% load: 170W (66%)
  - At 10% load: 121W (47%)
- Google data center
  - Mostly operates at 10% – 50% load
  - At 100% load less than 1% of the time
- Consider designing processors to make power proportional to load

# Pitfall: MIPS as a Performance Metric

- MIPS: Millions of Instructions Per Second
  - Doesn't account for
    - Differences in ISAs between computers
    - Differences in complexity between instructions

$$MIPS = \frac{Instruction\ count}{Execution\ time \times 10^6}$$

$$= \frac{Instruction\ count}{\dfrac{Instruction\ count \times CPI}{Clock\ rate} \times 10^6} = \frac{Clock\ rate}{CPI \times 10^6}$$

- CPI varies between programs on a given CPU

Scanned with CamScanner

# Content

# Concluding Remarks

- Cost/performance is improving
  - Due to underlying technology development
- Execution time: the best performance measure
- Power is a limiting factor
  - Use parallelism to improve performance

Scanned with CamScanner

# Chapter 4

## The Processor

*Adapted by Prof. Gheith Abandah*

# Contents

# Contents

4.6 Pipelined Datapath and Control (Review)

Five-Stage Pipeline

Pipeline Control

Pipeline Hazards

# Five-Stage Pipeline

**F**: Fetch instruction from the instruction ← inst الـ نجيب
memory ← inst. memory من

**D**: Decode instruction and read operands ← inst الـ تحليل
وقراءتها
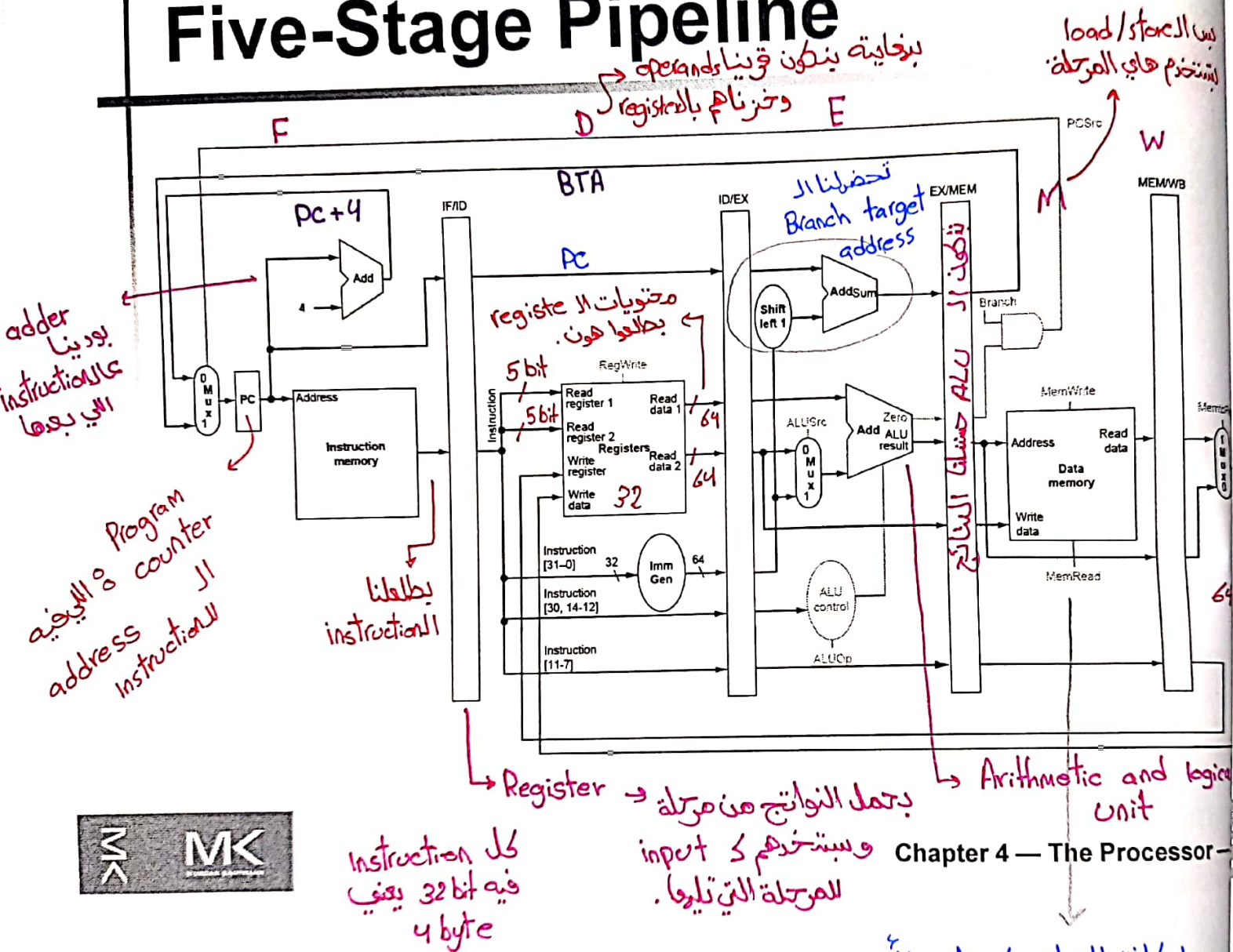
**E**: Execute operation or calculate address ← الحساب

**M**: Memory access ← اذا كانت, Load /
Store

**W**: Write result to the register

Scanned with CamScanner

# Five-Stage Pipeline

بنهاية بيكون قينا operands وخزناهم بالـ register

load/store نس الدنى يستخدم هاي المرحلة

F    D    E    W

BTA

PC + 4

IF/ID

Add

4

PC

تحضيرنا الـ Branch target address

ID/EX    EX/MEM    MEM/WB

AddSum

Shift left 1

Branch

adder بودينا عالـ instructions اللي بدها

Program counter الـ البيوخذ address instruction

Address

Instruction memory

بطلعنا الـ instruction

Instruction

RegWrite

5 bit

5 bit

Read register 1    Read data 1    64

Read register 2    Registers
Write    Read
register    data 2    64

Write
data    32

محتويات الـ registe بطالعوا هون.

ALUSrc

Add    Zero
ALU
result

M
U
X
1

0

MemWrite

Address    Read data

Data memory

Write data

MemRead

Memor

64

Instruction [31–0]    32    Imm Gen    64

Instruction [30, 14-12]    ALU control

Instruction [11-7]    ALUOp

جشنا الـ ALU بيتنفذ بنفذ

Arithmetic and logica unit

Register ← كل instruction فيه 32 bit يعني 4 byte

بتحمل النواتج من مرحلة وبستخدمهم كـ input للمرحلة التي تليها.

## Pipelined Control

لوكانت العملية load بنقرأ
لوكانت العملية store بنكتب
وهاي الـ data بتيجينا من Read data 2
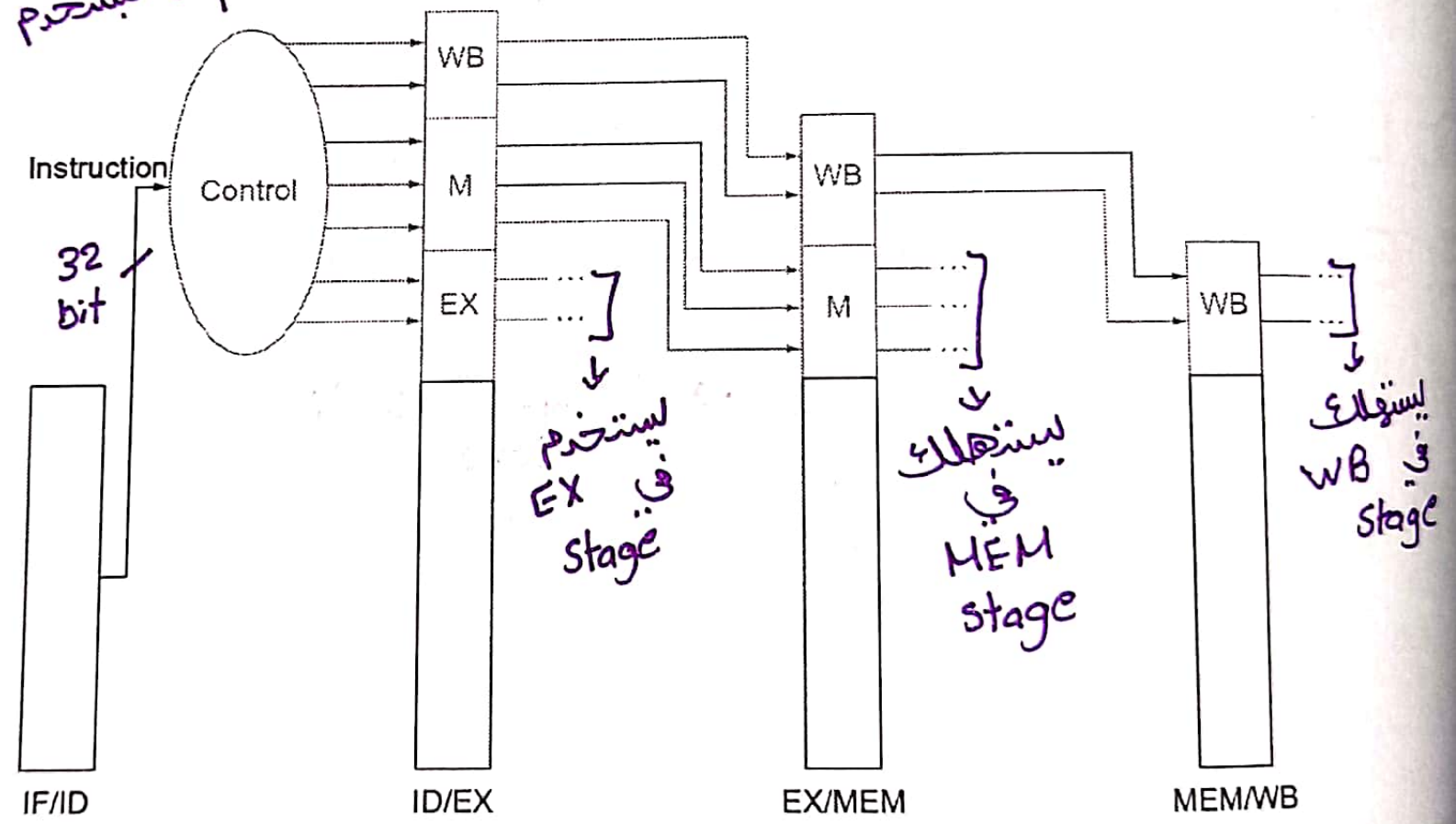
Control signals derived from instruction

# Control signals derived from instruction

- As in single-cycle implementation



opcode مع اللي هو 7-bit يستخدم
control لنطلع الـ
signals

Instruction

Control

32 bit

| WB |
| M |
| EX |

IF/ID

ID/EX

يستخدم في EX: Stage

EX/MEM

يستهلك في MEM stage

| WB |
| M |

| WB |

MEM/WB

يستغلك في WB stage

# Pipelined Control

*Regwrite: بتكون 1 عشان بدنا نكتب على الـ registers و 0 لما مابدنا نكتب مثل: الـ Branch والـ instruction مثل الـ store كمان

مثل الـ load و add/sub

*Branch: بتكون active لما يكون نوع الـ instr... branch

*MemtoReg: بتكون 0 لما الـ out يكون من الـ ALU و بتكون 1 لما الداتا بناخذها من الـ Mem مثل الـ load

*ALUsrc: لو كانت 0 خلا الـ ALU بجوا الداتا من الـ register و لو كانت 1 بدخل الـ Imm instruction

Chapter 4 — The Processor — 7

*ALUop: عبارة عن 2 bits بتحدد نوع العملية اللي تنفذها الـ ALU

# Hazards

*Mem Write: بتكون 1 لما بدنا نكتب على الـ Mem مثل الـ instruction store

*Mem read: بتكون 1 لما بدنا نقرأ من الـ Mem مثل الـ instruction load

- Situations that prevent starting the next instruction in the next cycle

*الفائدة الأساسية من الـ Pipelining: لتحسين الأداء Performance

1 Structure hazards → في مشكلة بالـ organisation → توقف خط العمل
  - A required resource is busy ← بنزيد من تصميمنا التحسينات

2 Data hazard ← تأخر اللي بتقرأ الـ data بينما ما تخلص اللي بتكتب
  - Need to wait for previous instruction to complete its data read/write

3 Control hazard → لها علاقة بالـ jump والـ Branch
  - Deciding on control action depends on previous instruction

Chapter 4 — The Processor — 8

Scanned with CamScanner

# Contents

MK

# Contents

## 4.7 Data Hazards: Forwarding versus Stalling

Hardware Solutions
- Data Hazards in ALU Instructions
- Load-Use Data Hazard

Software Solutions
- Code Scheduling

# Data Hazards in ALU Instructions

Consider this sequence:  *register instructions كل* 
*store الـ معدا*

```
sub   x2, x1,x3
and   x12,x2,x5        read from register and      store *
or    x13,x6,x2        write to memory
add   x14,x2,x2        read from memory and write : load *
sd    x15,100(x2)                 to register
      offset  ↵         ↳ base register
```
*كي نبقرأ من الـ وعي محتويات*
*x15 وبنا نخزنها بالـ Mem*

There are multiple true data dependencies read-after-write (RAW), on register x2.

*والـ Mem مكانه = محتويان 100 + x2*

We can resolve hazards with stalls or forwarding.

*Stalls 1- # حل الـ Hazards بطريقتين في*
*forward 2-*

*ولكن في hardware زيادة ر والحل الأفضل هو الـ Forward لأنه بزيد الـ Performence ويقلل الـ Stalls*
*بساحنا كبير transistors فمش مشكلة*

# Dependencies & Forwarding



| | Time (in clock cycles) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | CC 1 | CC 2 | CC 3 | CC 4 | CC 5 | CC 6 | CC 7 | CC 8 | CC 9 |
| Value of register x2: | 10 | 10 | 10 | 10 | 10/–20 | –20 | –20 | –20 | –20 |

Program execution order (in instructions)

*↳ بنكون الـ data الجديدة موجودة جيم*

sub x2, x1, x3

and x12, x2, x5

or x13, x6, x2

add x14, x2, x2

sd x15, 100(X2)

# Forwarding Paths

ID/EX    EX/MEM    MEM/WB

Registers

RF
W
M
M
u
x
ForwardA

ALU

RF
W
M
M
u
x
ForwardB

Rs1
Rs2
Rd

EX/MEM.RegisterRd

Data
memory

Forwarding
unit

MEM/WB.RegisterRd

**Chapter 4 — The Process**

بها تكتشف هل في data hazards
وهل في حاجة لل forwarding ولا
في حاجة بها تعطي control signal
ملائمة لل MUX لحتى ناخذ ال data من المكان الصحيح

# Load-Use Data Hazard

## Can't always avoid stalls by forwarding
- If value not computed when needed
- Can't forward backward in time!

Program
execution
order

| | 200 | 400 | 600 | 800 | 1000 | 1200 | 1400 |

## Slide 8 :

|     | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|-----|---|---|---|---|---|---|---|---|---|----|
| I1  | F | D | E | M | W |   |   |   |   |    |
| I2  |   | F | D | E | M | W |   |   |   |    |
| I3  |   |   | F | D | E | M | W |   |   |    |
| I4  |   |   |   | F | D | E | M | W |   |    |
| I5  |   |   |   |   | F | D | E | M | W |    |

*تسمى : "Multicycles Diagram"
لأنه احنا بنتابع بأكثر
من cycle وشوي مع
يحبين الـ instructions

↳ 5 Instructions

* لو مافي Stalls بنحتاج إلى 9 cycles لتنفيذ الـ instructions = 5 + 4 = 9

↳ مجموع عدد الـ cycles الزايدة بكل Instruction

↳ عدد الـ cycles اللي احتجهم لتنفيذ I1

* لو ما كان عندي pipeline كان احتجت 5 × 5 = 25 cycles لتنفيذ الـ instructions

## Slide 11 : <span style="color:red">Assume no forwarding (except through the Register file) and hazards are solved by stalls</span>

|                   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|-------------------|---|---|---|---|---|---|---|---|---|----|
| Sub $X_2, X_1, X_3$ | F | D | E | M | W |   |   |   |   |    |
| and $X_{12}, X_2, X_5$ |   | F | D | D | D | E | M | W |   |    |
| or $X_{13}, X_6, X_2$ |   |   | F | F | F | D | E | M | W |   |    |
| add $X_{14}, X_2, X_7$ |   |   |   |   |   | F | D | E | M | W  |
| sd $X_{15}, 100(X_2)$ |   |   |   |   |   |   | F | D | E | M  |

عملنا stall 1 → 2 cycles
بين ما احملنا القيمة
الجديدة
بنفضل بالـ F بين ما انخلص
مرحلة الـ D اللي قبلها

* بالـ stalls احتجت = 11 cycles = 4 + IC + stalls
كل instruction بدها 1 cycle لتعمل الـ W / عدد الـ cycles الزيادة بـ
(ثابت) 4 cycles احتون نعبي الـ pipeline و

* هون في 2 stalls بس لأنه الـ stall اللي عند الـ F لسبب نفس لسبب الـ Stall اللي
صار عند الـ D

* هاد الحل هو كويس لأنه نزّل من الـ performance

# Slide 12 : forwarding Solution :

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| sub $X_2$ , $X_1$ , $X_3$ | F | D | E | M | W | | | | | |
| and $X_{12}$ , $X_2$ , $X_5$ | | F | D | E | M | W | | | | |
| or $X_{13}$ , $X_6$ , $X_2$ | | | F | D | E | M | W | | | |
| add $X_{14}$ , $X_2$ , $X_2$ | | | | F | D | E | M | W | | |
| sd $X_{15}$ , $100(X_2)$ | | | | | F | D | E | M | W | |

\* خلصنا هم : 9 cycles \*

# Homework 1

A multicore process has a clock rate of 3.0 GHz and voltage of 1.0 v. Assume that, on average, it consumes 75 w of static power and 30 watt of dynami power when all its six cores are active. it executes a program consisting of $10^{10}$ instructions on a single core with the following instruction mix.

| Instruction type | CPI | Frequency |
|---|---|---|
| Arithmetic | 1.0 | 50% |
| Load/store → أبطا اشي | (2.0) | 40% |
| Branch instructions | 1.5 | 10% |

(a) Find the total execution time for this program on one core?

CPU time = IC × $\Sigma$(CPI × frequency) ÷ clock rate = $10^{10}$ × (1.0×0.5 + 2×0.4 + 1.5×0.1)

÷ $3×10^9$ = 4.833 Second

(b) when this program is parallized and run on the six cores, the number of Arithmetic and Load/store instructions per core is devided by 4, but the number of branches remains the same?

what is the execution time of the parallel program. = IC × $\Sigma$(CPI × frequency) ÷ clock rate

$10^{10}$ × (1.0×0.5 ÷4 + 2.0×0.4÷4 + 1.5×0.1) ÷ $3×10^9$ = 1.583 Sec

ممكن يوصل أحستم
منهيا يعني
6 مثلا

(c) what is the speedup of the parallel program? Speedup = 4.833 ÷ 1.583 = 3.05

(d) Assumes that the power consumed is proportional to the number of active cores. what are the energies consumeb by the serial program on the single core and the parallel program on the six cores? Energy of the serial program = exec.time × Power = 4.833 × ((75+30)÷6) = 85 Joules

Perfor... الكفاءة — Energy of the parallel program = exec.time × Power = 1.583 × (75+30) = 166 Joules
بتزداد

(e) what is the average capacitive load of this processor when the six cores are active?

Power = $\frac{1}{2}$ × Capacitive load × voltage$^2$ × frequency

30 = $\frac{1}{2}$ × capacitive load × $1^2$ × $3×10^9$ → capacitive load = 20 nano frades

(f) what is the average current drawn by this processor when the six cores are active?

power = current × voltage ⇒ 75+30 = current × 1 → current = 105 Ampere. → يعتبر كبير

ـ بتخلوها عن طريق تثرات ال ـ pins in parallel

(c) * لو مصو higher performence ⇒ استهلاك طاقة أكبر. *

# Load-Use Hazard Detection

- Check when using instruction is decoded in ID stage
- ALU operand register numbers in ID stage are given by
  - IF/ID.RegisterRs1, IF/ID.RegisterRs2
- Load-use hazard when شروط استخدام load-use
  - ID/EX.MemRead and  ↙  5 bits
  - 5 bits  ((ID/EX.RegisterRd = IF/ID.RegisterRs1) or (ID/EX.RegisterRd = IF/ID.RegisterRs2))
- If detected, stall and insert bubble

$D.RS_1 = E.Rd$
$D.RS_2 = E.Rd$     Stall

$E.MemRead$

Chapter 4 — The Processor — 15

# How to Stall the Pipeline

- Force control values in ID/EX register to 0
  - EX, MEM and WB do nop (no-operation)
- Prevent update of PC and IF/ID register
  - Using instruction is decoded again
  - Following instruction is fetched again
  - 1-cycle stall allows MEM to read data for ld
    - Can subsequently forward to EX stage

# Load-Use Data Hazard

Time (in clock cycles)

CC 1   CC 2   CC 3   CC 4   CC 5   CC 6   CC 7   CC 8   CC 9

Program
execution
order
(in instructions)

ld x2, 20(x1)

and becomes nop

and x4, x2, x5

or x8, x2, x6

add x9, x4, x2

bubble

Stall inserted here

# Datapath with Hazard Detection

Hazard detection unit

ID/EX.MemRead

Rs1

Rs2

E.Rd

IF/IDWrite

PCWrite

PC

Instruction memory

IF/ID

Instruction

Control

Registers

Mux

ID/EX

WB

M

EX

Mux

ForwardA

ALU

Mux

ForwardB

EX/MEM

WB

M

Data memory

MEM/WB

WB

IF/ID.RegisterRs1
IF/ID.RegisterRs2
IF/ID.RegisterRd

Rd

Rs1
Rs2

Forwarding unit

# Stalls and Performance

- Stalls reduce performance
  « بنجتاجه بحالات »
  - But are required to get correct results
- Compiler can arrange code to avoid hazards and stalls
  - Requires knowledge of the pipeline structure

→ slide 16

\* شو بنحتاج نعمل بال Stall cycle :
disable نقتل PC الا نزوح Signal -1
disable    IR → Rigister after PC -2
No-operation (0) to Execution -3

# Code Scheduling to Avoid Stalls

- Reorder code to avoid use of load result in the next instruction
- C code for a = b + e; c = b + f;

↓ RISC-v

```
        ld    x1, 0(x0)      b        ld    x1, 0(x0)
        ld    x2, 8(x0)      e        ld    x2, 8(x0)
stall   add   x3, x1, x2              ld    x4, 16(x0)
Statement1 sd x3, 24(x0) a            add   x3, x1, x2
        ld    x4, 16(x0) f            sd    x3, 24(x0)
stall   add   x5, x1, x4              add   x5, x1, x4
Statment 2 sd x5, 32(x0) c            sd    x5, 32(x0)
```

13 cycles          11 cycles

4+7+2=13           4+7= 11

# Contents
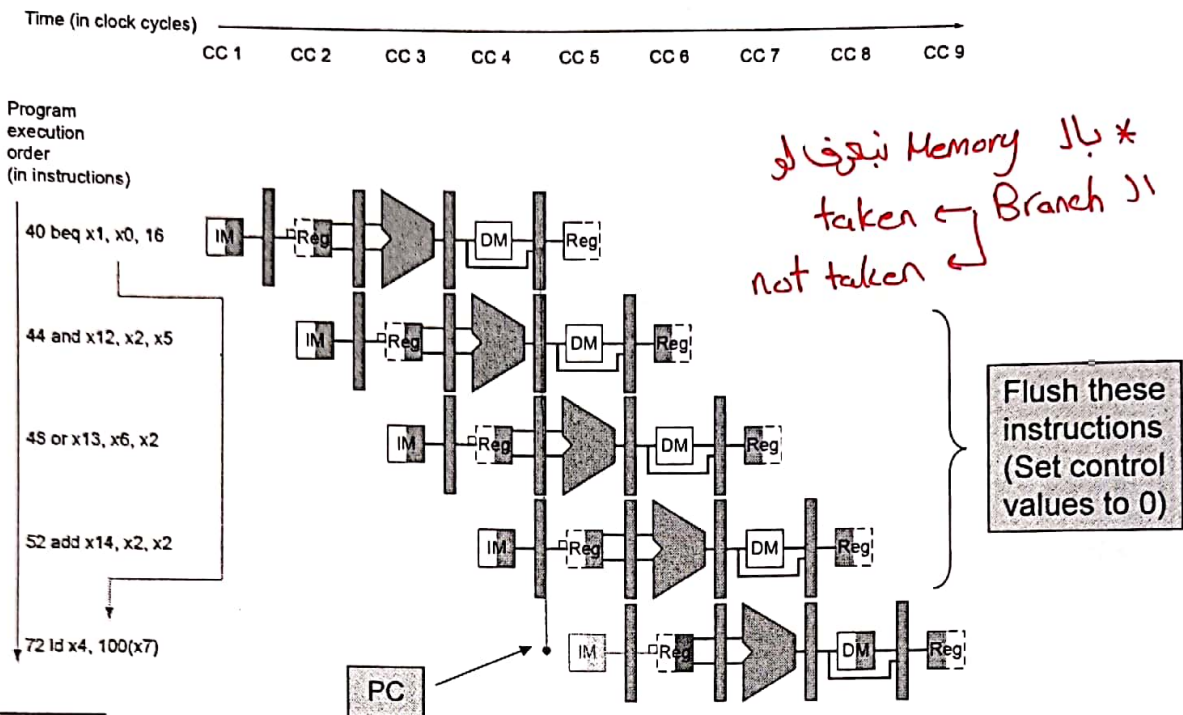
ЗМ MK

# Contents

لعملو ا delay فكيف وقلليم stalls عن طريق

4.8 Control Hazards

Branch Hazards

Reducing Branch Delay

Branch Prediction

Dynamic Branch Prediction

Calculating Branch Target

Imprecise Exceptions

# Branch Hazards

- If branch outcome determined in MEM

Time (in clock cycles)

CC 1   CC 2   CC 3   CC 4   CC 5   CC 6   CC 7   CC 8   CC 9

Program
execution
order
(in instructions)

40 beq x1, x0, 16

44 and x12, x2, x5

48 or x13, x6, x2

52 add x14, x2, x2

72 ld x4, 100(x7)

PC

* بالنسبة لل Memory نتعرف على
taken ⟋ البranch ال
not taken ⟍

Flush these
instructions
(Set control
values to 0)

# Reducing Branch Delay

- Move hardware to determine outcome to ID stage
  - Target address adder
  - Register comparator
- Example: branch taken

```
36:  sub  x10, x4, x8
40:  beq  x1,  x3, 16   // PC-relative branch
                        // to 40+16*2=72
44:  and  x12, x2, x5
48:  orr  x13, x2, x6
52:  add  x14, x4, x2
56:  sub  x15, x6, x7
     ...
72:  ld   x4, 50(x7)
```

# Example: Branch Taken



and x12, x2, x5        beq x1, x3, 16        sub x10, x4, x8        before<1>

Clock 3

# Example: Branch Taken



ld x4, 50(x7)        Bubble (nop)        beq x1, x3, 16        sub x10, ...        before<

Clock 4

# Branch Prediction

- Longer pipelines can't readily determine branch outcome early
  - Stall penalty becomes unacceptable
- Predict outcome of branch
  - Only stall if prediction is wrong
- In RISC-V pipeline
  - Can predict branches not taken *اقترض الطرق أسهل* *branches not taken*
  - Fetch instruction after branch, with no delay

---

*Branch Delay كبير* *طويل فحيطلع ال Pipline لو كان ال Predict not taken لأنه طريقة*

# More-Realistic Branch Prediction

1) Static branch prediction → *Static بينصرف* *يعني ال Branch Prediction unit*
   - Based on typical branch behavior *مش بناء على الحالة*
   - Example: loop and if-statement branches *الأصلية اللي هتردف تتعلم منها.*
     - Predict backward branches taken
     - Predict forward branches not taken

2) Dynamic branch prediction
   - Hardware measures actual branch behavior ↳ *عكس ال Static، بنتعلم من ال actual اللي بحصل*
     - e.g., record recent history of each branch
   - Assume future behavior will continue the trend *بنستفيد من الماضي عشان نتنبأ في الحاضر*
     - When wrong, stall while re-fetching, and update history

# Dynamic Branch Prediction

- In deeper and superscalar pipelines, branch penalty is more significant
- Use dynamic prediction
  - Branch prediction buffer (aka branch history table)
  - Indexed by recent branch instruction addresses
  - Stores outcome (taken/not taken)
  - To execute a branch
    - Check table, expect the same outcome
    - Start fetching from fall-through or target
    - If wrong, flush pipeline and flip prediction

MK

أبسط
dynamic Branch طريقة بالـ
prediction

# Branch History Table (BHT)

NT: 0
T: 1

← Branch prediction هذه يعطيني
bits

## One-Level Branch Predictor   Memory



64 bit
Branch Address

13 bit بالطبيعي

k bits

n-bit counter

| |
|---|
| 0 |
| 1 |
| 0 |
| × × |
| |
| ... |
| |

الناتج
1 و 0

prediction bits

Table size = $n \times 2^k$ bits

كمثال ← $2 \times 2^{10} = 2k$ bits.

MK

# 1-Bit Predictor: Shortcoming

- Inner loop branches mispredicted twice!

```
outer: ...
       ...
inner: ...
       ...
       beq ..., ..., inner
       ...
       beq ..., ..., outer
```

- Mispredict as taken on last iteration of inner loop
- Then mispredict as not taken on first iteration of inner loop next time around

# 2-Bit Predictor

- Only change prediction on two successive mispredictions

# Calculating the Branch Target

- Even with predictor, still need to calculate the target address
  - 1-cycle penalty for a taken branch
- Branch target buffer
  - Cache of target addresses
  - Indexed by PC when instruction fetched
    - If hit and instruction is branch predicted taken, can fetch target immediately

عشان بالمستقبل لو مرت على بنفس ال branch وكان ال address تاعها نفس
الي مخزنك عندي فبوديني ال BTA عن طريق ال MUX

كل ما يمر على Branch بخزن ال address تبعها ولدين بترج
من ال PC في ال Address هاد ال Address بستخدمه لنوجد BTB

# Branch Target Buffer (BTB)

كبارق عن ال Addresses في ال Memory
ال branches اللي مروا على وفيه BTA
إنه هدول ال branches اللي مروا على وبن راحوا



| | PC of instruction to fetch | |
| --- | --- | --- |
| | 64 ↓ Look up | Predicted PC |

Number of entries in branch-target buffer

64

64 ↓

64 → No: Instruction is not predicted to be branch; proceed normally

2bit
1 bit

Branch predicted taken or untaken

Yes: then instruction is branch and predicted PC should be used as the next PC

ممكن يخزنوا معاه ال BHT بس. بولي الأوقات ما بتكون موصولين
سوا غال با مفصولين

# Contents

بخلينا نوقف شغلناعلى البرنامج و نروح نحله

MK

# Contents

## 4.9 Exceptions

Exceptions and Interrupts

Handling Exceptions

Exceptions in a Pipeline

Exception Example

Multiple Exceptions

# Exceptions and Interrupts

- "Unexpected" events requiring change in flow of control
  - Different ISAs use the terms differently
- Exception ١)  ← يعني داخل ال pipline
  - Arises within the CPU  وهو بيتخذ ممكن يصير error بحتاج معالجة
    - e.g., undefined opcode, syscall, ... و division by zero
- Interrupt
  - From an external I/O controller  ← بيجونا من الخارج
- Dealing with them without sacrificing performance is hard

# Handling Exceptions

ال Instruction اللي عملت Exep. →

- Save PC of offending (or interrupted) instructio
  - In RISC-V: Supervisor Exception Program Counter (SEPC)

- Save indication of the problem
  - In RISC-V: Supervisor Exception Cause Register (SCAUSE)
  - 64 bits, but most bits unused
    - Exception code field: 2 for undefined opcode, 12 for hardw malfunction, ...
- Jump to handler
  - Assume at $0000\ 0000\ 1C09\ 0000_{hex}$

# An Alternate Mechanism

Vectored Interrupts
- Handler address determined by the cause

Exception vector address to be added to a vector table base register:
- Undefined opcode $\qquad$ 00 0100 0000$_{two}$
- Hardware malfunction: $\quad$ 01 1000 0000$_{two}$
- ...: $\qquad\qquad\qquad\qquad$ ...

Instructions either
- Deal with the interrupt, or
- Jump to real handler

MK

# Handler Actions

Read cause, and transfer to relevant handler

Determine action required

If restartable
- Take corrective action
- use SEPC to return to program

Otherwise
- Terminate program
- Report error using SEPC, SCAUSE, ...

MK

Slide 148 Load-use data hazard

ثابت ← Instruction stalls

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | → 4 + 2 + 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ld X₁ , 0(X₂) | F | D | E | M | W | | | | | | = 7 cycles |
| sub X4, X₁ , X5 | | F | D | D | E | M | W | | | | |

ld $X_1$ , 0($X_2$)

sub $X_4, X_1, X_5$

\* الـ Sub والـ add بتجهز الـ data بالـ Execute stage       لحمانا 1 stall لشان الـ data بالـ load

\* الـ load بتجهز الـ data بالـ Memory بنهايتها       بنجهز نهاية الـ Memory ، بعين عملنا

Forwarding لشان نحط القيم الجديدة .

---

5 bit     5 bit       \* لحل الـ Load-use بنحتاج نعمل 1 stall و Forwarding .

D.RS₂ = E.Rd   أو   D.RS₁ = E.Rd       \* الدائرة اللي بنحتاجها لشوفين :

5-bits ← A₄ A₃ A₂ A₁ A₀       → A = B او بتطلعي 1



B₄ B₃ B₂ B₁ B₀

Slide 20: Code Scheduling to Avoid Stalls:

EX:
```
ld   X1, 0(X2)
add  X3, X1, X4
ld   X5, 0(X6)
```
↳ independent

إعادة ترتيب →

```
ld   X1, 0(X2)    F D E M W
ld   X5, 0(X6)    F D E M W
add  X3, X1, X4     F D E M W
```

↓

هون بهياء ترتيب مابحتاج أعمل

Stall بس Forwarding

---

Slide 23: Branch Hazards

Solving branches in the Memory stage

* منحولهم لـ null instruction

لإنه مابدنا يتنفذوا

* بالـ Memory

بنعرف إذا الـ Branch

taken أو not taken

بيجايتها

* افترضنا أنه

الـ Branch is taken

| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 40 | beq X1, X0, 16 | F | D | E | M | W | | | | | | |
| 44 | and X12, X2, X5 | | F | D | E | n | ∧ | | | | | |
| 48 | or X13, X6, X2 | | | F | D | ∧ | ∧ | ∧ | | | | |
| 52 | add X14, X2, X2 | | | | F | ∧ | ∧ | n | ∧ | | | |
| 72 | ld X4, 100(X7) | - | | | | F | D | E | M | W | | |

↳ target address = 40 + 2×16 = 72

لأنه بنزيد/بنضرب في 2 half bytes

↳ في cycle5 وبداية cycle4 الـ بتوقف

لأن الـ Branch is taken فبعمل F

* هون الـ Performence بيقل لإنه كاننا عندي 3 stall cycles

Branch delay = 3

* ولتقليل الـ Branch delay → هنا لازم بأنه بدل بالـ Memory stage يجوز بالـ

بظلوا يجوز بالـ Decode stage

Scanned with CamScanner

Slide 24 :

Soluing branches in the Decode stage → صرنا هون نعرف نتيجة الـ Branch

Branch is taken

| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 40 | beq $X_1, X_0, 16$ | F | D | E | M | W | | | | | |
| 44 | and $X_{12}, X_2, X_5$ | | F | ∩ | ∩ | ∩ | ∩ | | | | |
| 48 | or $X_{13}, X_6, X_2$ | | | | | | | | | | |
| 52 | add $X_{14}, X_2, X_2$ | | | | | | | | | | |
| 72 | ld $X_4, 100(X_7)$ | | | F | D | E | M | W | | | |

Branch delay = 1   *

* التحسين الجديد هو أن أقوم بعمل → Branch prediction
بخلي الـ Hardware يتنبأ بالنتيجة.

Slide 28 : predict Not taken

Soluing branches in the Decode stage

① Assume branch is not taken

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| beq $X_1, X_0, L$ | F | D | E | M | W | | | | | | → توقفنا وبلاش |
| $I_2$ | | F | D | E | M | W | | | | | وكان توقفنا صح |
| | | | | | | | | | | | وما ضيعنا أي |
| | | | | | | | | | | | cycle |
| L | $I_T$ | | | | | | | | | | → Branch Delay |
| | | | | | | | | | | | = 0 |

② Assume branch is taken          avg Branch Delay = $\frac{1}{2}$ cycle

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| beq $X_1, X_0, L$ | F | D | E | M | W | | | | | | → Branch |
| $I_2$ | | F | ∩ | ∩ | ∩ | ∩ | | | | | Delay = 1 |
| | | | | | | | | | | | |
| L | $I_T$ | | | F | D | E | M | W | | | |

## Slide 28 : Branch Predection Unit

F

IF/ID

BTA

Pc+4

M U X

Pc

Instruction Memory

I

4 →

+

↑ T, NT

BTA

| Branch prediction unit | → ‎إضافته للتصميم وشغلته ← من الـ Address‏ |

‎الـ instruction اللي قاعدين بنجيبها بيتنبأ إذا‏
taken or not taken ← Branch ‎الـ‏
‎وأيضا مسؤولة عن إعطاء الـ BTA‏

---

## Slide 30 : BHT → n=1 / n=2   ↓ n=1
1-bit predictor

taken

not taken

Not take 0

taken 1

taken

Not taken

‎لو افترضنا الـ Prediction taken‏
‎وطلعت فعلاً taken بنخليها taken‏
(‎وهي جهاز للحالات الباطنية‏)

0 → 999

Assume 1000 iterations for the inner loop and for the outer loop

outer :... → 1000 ‎مرة‏

inner : ... 1000 ‎مرة‏
× 1000
...

...  1000000 ‎فتنفذ‏

bea, --, ---, inner

bea, --, ---, outer

‎لأنه المفروض يطلع‏
‎ما ننفذها‏
‎أخطأ‏
‎مرتين‏

| Iteration | 997 | 998 | 999 | 0 | 1 |
|---|---|---|---|---|---|
| Prediction Result | T/C | T/C | T/I | NT/I | T/C |

* T: taken    /  NT: Not taken

C: correct prediction / I: incorrect prediction

Scanned with CamScanner

* الـ one bit بغير رأيه بسرعة فاتجوننا, اكن الـ 2 bit يغير رأيه بطريقة أبطأ.

## Slide 32:
## 2-bit Predictor

Predict taken 10 — Not taken → Predict taken 11

taken ←

T 1 → most sig bit

Predict Not taken 00 (عميق) ← Not taken — Predict Not taken 01 (خفيف)

taken →

NT 0 → most sig bit

outer: .....
inner: ...

beq,--,--, inner

beq -,..., outer

| Iteration | 997 | 998 | 999 | 0 | 1 |
|---|---|---|---|---|---|
| Prediction result | T/c | T/c | T/1 | T/c | T/c |

1 miss Prediction

Slide 38: Handling Exceptions

* ميزة هاي الطريقة إنها simple hardware
بس ال Software أصعب

I1       I3 في Exception صار لنفرض *

I2

I3 *   و   15 و 14 ناكمل المفروض

I4       اللي Address خال لنزح بل

I5     Exception handling    ال عنده موجود

المخصص Address ال هاد RISC-V وبال

0000 0000 1c09 0000 hex

┌─────────────┐
│ Exception   │
│ Handling    │
│ Routine     │
└─────────────┘

برنامج جزء من
ال operating system
بقدر يتعامل مع
ال Exceptions

┌──────────────────┐
│ SEPC, SCAUSE     │  →   داخل ال Prossesor لازم نضيف Additional Register
└──────────────────┘
         SCAUSE , SEPC اللي هم

Exception ال صار عليها اللي ال instruction ال في يخزن بي لازم Exception ال نخزن في ال Address لازم بس هدول
وبنفس الوقت ببي أخزن سبب ونوع ال Exception

* Exception Handling Routine: خزناه اللي ال Exception سبب من يستفيد رح
ال instruction يرجع بيقدر ال SEPC ال من وبعين ال Exception ويعالجه ال SCAUSE بال
الأصلي البرنامج تنفيذ ويتابع عندها توقفنا اللي

* ال Exception Handling Routine دايما مش بقدر يرجع ال instruction اللي توقفنا عندها ،
إذا قدر يحل المشكلة وبنقدر نكمل البرنامج بيرجع ، وبعض المرات بروح مكان تاني بال operating system ولكان ما يقدر يحل المشكلة مثل مشكلة : code الوصول illegal
وساعتها ال OS بطلعلا user رسالة بتقوله انه عندك مشكلة ما بنقدر نكمل برنامجك

Slide 40:

I1          undefined opcode    00 0100 0000 two

I2          Hardware malfunction: 01 1000 0000 two

I3     * هاي الطريقة حسب نوع ال Exception بنروح ل Address مختلف .

I4     * ميزة هاي الطريقة بتخليه أسرع

I5     * ال Software أسول .

┌──────────────────┐
│ SEPC, SCAUSE     │
└──────────────────┘

# Exceptions in a Pipeline

- Another form of control hazard
- Consider malfunction on add in EX stage

    add x1, x2, x1

    - Prevent x1 from being clobbered
    - Complete previous instructions
    - Flush add and subsequent instructions
    - Set SEPC and SCAUSE register values
    - Transfer control to handler
- Similar to mispredicted branch
    - Use much of the same hardware

**MK**

# Pipeline with Exceptions



ال address
اللي بروحيله
لو حدث exceptions

# Exception Properties

■ Restartable exceptions → يعني بنقدر نرجع للـ instruction الي صار عليها Exception بعد ما نحله

- Pipeline can flush the instruction
- Handler executes, then returns to the instruction

  - Refetched and executed from scratch → عكسهم الـ unrestartable

■ PC saved in SEPC register

- Identifies causing instruction → الـ instruction اللي صار عندها exception

# Exception Example

■ Exception on add in

```
40      sub   x11, x2, x4
44      and   x12, x2, x5
48      orr   x13, x2, x6
4c      add   x1,  x2, x1
50      sub   x15, x6, x7
54      ld    x16, 100(x7)
...
```

هدول لازم يشفذوا ويخلصوا
" هون صار فيه exception "
هدول بنعملهم Flash

■ Handler

```
1C090000    sd   x26, 1000(x10)
1C090004    sd   x27, 1008(x10)
...
```

بالنهاية رح يعمل
ld x27
ld x26

بالعادة الـ Exception handling subroutine بيبلشوا بـ sd : لإنه حتى يشتغل بده reg

يقرأ ويكتب مفهوم وبنفس الوقت مابده يخرب الـ data اللي بالبرنامج الأصلي وما ينفع فبعملهم Save

Scanned with CamScanner

# Exception Example



ld x16, 100(x7)   sub x15, x6, x7   add x1, x2, x1   or x13, . . .   and x1?

Clock 6

# Exception Example

* الدكتور بعمل Flash عند W مش عند الـ M لإنه هيك أدق وأصح ليعطي مجال نخلص الـ Instruction اللي قبل

الدكتور بعمل Flash عند W مش عند الـ M لإنه هيك أدق وأصح ليعطي مجال نخلص الـ Instruction اللي قبل مش زي السلايد باد M بعملوا



sd x26, 1000(x0)   bubble (nop)   bubble   bubble   or x1?

Clock 7

Scanned with CamScanner

# Multiple Exceptions

- Pipelining overlaps multiple instructions
  - Could have multiple exceptions at once
- Simple approach: deal with exception from earliest instruction
  - Flush subsequent instructions
  - "Precise" exceptions → زمان ماكانوا يقدروا يعملوا هيك؟
- In complex pipelines
  - Multiple instructions issued per cycle
  - Out-of-order completion
  - Maintaining precise exceptions is difficult!

كانوا يعملوا هيك؟

# Imprecise Exceptions → Processors

كانت لما يصير exception
ما نعرف مين خلص ومين ما خلص وهو عملية معقدة و بعض المرات لو حل exception

- Just stop pipeline and save state
  ما بنقدر نرجع للبرنامج الأصلي
  لأنه ما بنعرف لو كل الـ instruction
  - Including exception cause(s) اللي قبل الـ exception خلصوا ولا لا
    و لا الـ اللي بعدها
- Let the handler work out
  - Which instruction(s) had exceptions
  - Which to complete or flush
    - May require "manual" completion
- Simplifies hardware, but more complex handler software
- Not feasible for complex multiple-issue out-of-order pipelines

# Contents

*Multiple نفذ instructions every cycle وهي تعني هو*

*modern بال Processors*

# Contents

## 4.10 Parallelism via Instructions

Instruction-Level Parallelism (ILP)

Multiple Issue

Static Multiple Issue

VLIW

Scheduling Static Multiple Issue

Loop Unrolling

Dynamic Multiple Issue

Register Renaming

Speculation

Why Do Dynamic Scheduling

*أفضل لأن الجدولة الديناميكية*

# Instruction-Level Parallelism (ILP)

Pipelining: executing multiple instructions in parallel → $Time = Ic \times CPI \times Tinst$

To increase ILP

① ■ Deeper pipeline T↓

الي أخذناه قبل انه
لنفذوا ب Multiple instructions
بمراحل مختلفة عشان وحدة لبس

   Less work per stage ⇒ shorter clock cycle

② ■ Multiple issue

   Replicate pipeline stages ⇒ multiple pipelines

   Start multiple instructions per clock cycle

   CPI < 1, so use Instructions Per Cycle (IPC)

   E.g., 4GHz 4-way multiple-issue → instruction per cycle

   ═ 16 BIPS, peak CPI = 0.25, peak IPC = 4

   But dependencies reduce this in practice

مثلا يكون

```
F D E
 F D E
  F D
   F D
    ¦
```
تحت بعض عادي

$f * IPC$

Time قلّ كلّما CPI قلّ

لـ ينزل الـ

بقرب الـ

4 inst:
```
F D E
F D E
F D E
F O E
```

* $CPI = \dfrac{1}{IPC} = \dfrac{1}{4} = 0.25$

# Multiple Issue

العبئ بيجي على الـ compiler أو المبرمج
① Static multiple issue → والمبرمج بحط الـ instruction اللي ما بعتمدوا عبضن اللي ممكن تنفيذهم مع بوض بمجموعات

■ Compiler groups instructions to be issued together والـ Hardware

■ Packages them into "issue slots" ياخذ هاد الـ issue وما بيحتاج slots

■ Compiler detects and avoids hazards يفحصه لأنه الـ Compiler detects and avoids hazards

② Dynamic multiple issue

■ CPU examines instruction stream and chooses instructions to issue each cycle أقوى و higher performance

■ Compiler can help by reordering instructions

■ CPU resolves hazards using advanced techniques at runtime

الـ CPU بنفسها بتجيب بتنجيب الـ instruction و بتعملها Decode و يعرف ايش الـ dependences
بنبهم Hardware و الـ يرضها ← chooses instructions to issue each cycle

# Static Multiple Issue

- Compiler groups instructions into "issue packets" → <span dir="rtl">بجمع</span> issue Packets ← independent instructions

  - Group of instructions that can be issued on a single cycle
  - Determined by pipeline resources required

- Think of an issue packet as a very long instruction

  - Specifies multiple concurrent operations
  - ⇒ Very Long Instruction Word (VLIW)

<span dir="rtl">مثال Itanium</span>  <span dir="rtl">هو ال ما بشبه ال</span>  Static Multiple issue

MK

Ex: add X₁, X₂, X₃
$\left.\begin{array}{l} \text{add } X_1, X_2, X_3 \\ \text{Sub } X_4, X_1, X_5 \\ \text{add } X_6, X_7, X_8 \end{array}\right]$ → add X₁, X₂, X₃ / add X₆, X₇, X₈ → issue Packets

dependent → <span dir="rtl">ما بنقدر ننفذهم مع بعض</span>
→ independent <span dir="rtl">لهذا أنخذهم مع بعض وأحطهم في</span> issue packet

## VILW

<span dir="rtl">ممكن نعتبره نوع من</span> <span dir="rtl">الـ</span> Static Multiple issue <span dir="rtl">الـ</span> VLIW →

<span dir="rtl">الوحدة أكبر دكيش</span> instruction <span dir="rtl">الـ</span>

(very long instruction word, 1024 bits!) <span dir="rtl">من</span> 32 bit

<span dir="rtl">وهاي الـ</span> Instruction

Specify multiple concurrent operations

Processor <span dir="rtl">ال</span> <span dir="rtl">يعني ننقول الـ</span> one operation <span dir="rtl">اعل أكثر منها</span> <span dir="rtl">بنفس الوقت</span> (concurrent)

VLI instruction

Cache/ memory <span dir="rtl">فيها</span> → Fetch unit <span dir="rtl">بنجيب الـ</span> VLI

Single multi-operation instruction

Multi-operation instruction ← <span dir="rtl">يعني</span> instruction <span dir="rtl">فيها</span> Multiple operations <span dir="rtl">بكون</span> 128 or 256 bit

EU ALU | EU FP | ----- | EU MeM

Register file

VLIW approach

Compiler <span dir="rtl">الـ</span> <span dir="rtl">هو ايش</span> operations <span dir="rtl">الـ</span> <span dir="rtl">اللي ممكن ننفذهم مع بعض وحطهم</span> <span dir="rtl">بـ</span> VLI <span dir="rtl">وكل</span> operation <span dir="rtl">بشغد على</span> Execution unit <span dir="rtl">خاصة</span>

MK

# Scheduling Static Multiple Issue

## Compiler must remove some/all hazards

- Reorder instructions into issue packets
- No dependencies with in a packet
- Possibly some dependencies between packets → instruction يكون فيه Packet والـ Packet ممكن بين الـ
  بتعتمد على بعضها
  - Varies between ISAs; compiler must know!
- Pad with nop if necessary → independent instruction إذا ما قدر يلاقي

كافية يحطوهم به Packet واحد يعبي باقي الـ Packet + no operations

يعني وجي المحلات الفاضية

# RISC-V with Static Dual Issue

issue packet يعني الـ

- Two-issue packets — 2 instructions بكون فيه
  - One ALU/branch instruction — load/store ووحدة ALU/ وحدة بتكون / branch
  - One load/store instruction — ( nop ) ولو ما قدر يلاقي بخفيف
  - 64-bit aligned
    - ALU/branch, then load/store
    - Pad an unused instruction with nop

مفترض. إنه ما في Stalls ولا في hazards

PC+8

| Address | Instruction type | Pipeline Stages | | | | | | |
|---------|------------------|-----|-----|-----|-----|-----|-----|-----|
| n | ALU/branch | IF | ID | EX | MEM | WB | | |
| n + 4 | Load/store | IF | ID | EX | MEM | WB | | |
| n + 8 | ALU/branch | | IF | ID | EX | MEM | WB | |
| n + 12 | Load/store | | IF | ID | EX | MEM | WB | |
| n + 16 | ALU/branch | | | IF | ID | EX | MEM | WB |
| n + 20 | Load/store | | | IF | ID | EX | MEM | WB |

issue Packet
issue Packet
issue Packet

بصير الـ throughput أفضل لإنه بننفذ 2 في cycle

2 IPC → instruction per cycle

$\frac{1}{2}$ CPI → clock per Instruction

# RISC-V with Static Dual Issue

بينفذ ال add والا sub والا Branch

4 read port بربعين
2 write port و

Memory address
Rs + Imm

# Hazards in the Dual-Issue RISC-

- More instructions executing in parallel
- EX data hazard → من ال issue و ال اللي بعديه issue ممكن يكون فيه dependent ونحتاج Forwarding
  - Forwarding avoided stalls with single-issue
  - Now can't use ALU result in load/store in same packe

```
add    x10, x0, x1
ld     x2, 0(x10)
```
مابنقدر نحطهم مع بعض لأنه بيتفذوا على بعض
  - Split into two packets, effectively a stall
- Load-use hazard
  - Still one cycle use latency, but now two instructions
- More aggressive scheduling required

# Forwarding in Dual-Issue RISC-V

In addition to forwarding from M and W to E, there are additional forwarding paths among the two pipelines, e.g.:

- From W in memory pipeline to E in ALU pipeline

```
ld    x31, 0(x20)
add   x31, x31, x21
```

أخذناه قبل →

- From M in ALU pipeline to M in memory pipeline

```
add   x31, x31, x21
sd    x31, 0(x20)
```

---

# Scheduling Example

* حصلنا على IPC = 1.25 عشان الnop
فلو حطينا بدل الnop → instructions
بزيد الـ IPC فلذلك بنستخدم حل أحسن واللي هو loop unrolling

- Schedule this for dual-issue RISC-V

```
Loop: ld    x31,0(x20)     // x31=array element
      add   x31,x31,x21    // add scalar in x21
      sd    x31,0(x20)     // store result
      addi  x20,x20,-8     // decrement pointer
      blt   x22,x20,Loop   // branch if x22 < x20
```

2 issue Packets

ممكن نشط بـ 1 برضه ← و حطيناها أول وحدة لأننا أول عملية

|       | ALU/branch          | Load/store      | cycle |
|-------|---------------------|-----------------|-------|
| Loop: | nop                 | ld   x31,0(x20) | 1     |
|       | addi  x20,x20,-8    | nop             | 2     |
|       | add   x31,x31,x21   | nop             | 3     |
|       | blt   x22,x20,Loop  | sd   x31,8(x20) | 4     |

load use hazard
ما قررنا بالثانية لأنه

- IPC = 5/4 = 1.25 (c.f. peak IPC = 2)

instruction في
issue Packet

الـ addi قبلها بصل
فلازم أخليها هو 0 عشان
نعكس أثر الـ addi

بس ممكن أحطها بـ 3 برضه
لأنها بتيجي بعد الـ add

Scanned with CamScanner

# Loop Unrolling

- Replicate loop body to expose more parallelism
- فوائده ■ Reduces loop-control overhead
- Use different registers per replication
  - ■ Called "register renaming" → بخلينا نتقلب على أنواع
    dependence الـ من إضافية
  - محتاج لعمل ■ Avoid loop-carried "anti-dependencies"
    loop unrolling
    - Store followed by a load of the same register
    - Aka "name dependence", write-after-read
    - Or "output dependence", write-after-write
      - ■ Reuse of a register name

MK

if ( -- ─, ) i+=1
a[i] = a[i]+3;

مثال :

i+=2
if ( -- --i ) {
a[i] = a[i]+3;
a[i+1] = a[i+1]+3;
}

instruction : بطريقة أحسن وبكبرال وصا ويزيد يعمله

# Unrolling Steps

بعملها المبرمج و
Compiler الـ بعملها الأيام هاي

أنت
بتختار م

1. Replicate the loop instructions n times
2. Remove unneeded loop overhead
3. Modify instructions
4. Rename registers
5. Schedule instructions

# Loop Unrolling Example

|        | ALU/branch            | Load/store          | cycle |
|--------|-----------------------|---------------------|-------|
| Loop:  | addi x20,x20,-32      | ld   x28, 0(x20)    | 1     |
|        | nop                   | ld   x29, 24(x20)   | 2     |
|        | add x28,x28,x21       | ld   x30, 16(x20)   | 3     |
|        | add x29,x29,x21       | ld   x31, 8(x20)    | 4     |
|        | add x30,x30,x21       | sd   x28, 32(x20)   | 5     |
|        | add x31,x31,x21       | sd   x29, 24(x20)   | 6     |
|        | nop                   | sd   x30, 16(x20)   | 7     |
|        | blt x22,x20,Loop      | sd   x31, 8(x20)    | 8     |

32 general

32 Floating

register بس عادي لإنه هنا كثير

IPC = 14/8 = 1.75

ماوصلنا 2 عشان لسا في nop

بس هاي طريقة: أحسن

X28,X29,X30,X31,X21 لأنه قبل كنا نستخدم بس ...,X20,X31 بس هلا هنا نستخدم بالإضافة السابق

Closer to 2, but at cost of registers and code size

خلو يهمني الـ Performence مايهمني لو حجم البرنامج كبر شوي

هو يحسن الـ Performence بس على حساب الـ registers والـ
code size ولكن هاي المشاكل مقبولة والـ Mem هاي الأيام رخيصة وكبيرة

يعني استهلكنا Memory usage

MK

32 bit

كبر حجمه أكثر من 3 أضعاف

32 bit

\* الـ code اللي قبل : 5 x 4 = 20 byte instructions

issue Packets

8 x 2 x 4 = 64 bytes : الـ code الجديد

instruction per issue packets

## Slide 41 : Exceptions in a pipeline :

بيي أفترض : إنه بس اعلوا E في مشكلة exception إي جمعت int + int طلع كبير ما بقدر أخزنه بـ $X_1$

| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | I1 | F | D | E | M | W | | | | | | | | |
| | add $X_1$, $X2$, $X_1$ | | F | D | Ë | M | n | | | | | | | |
| | I3 | | | F | D | E | n | | | | | | | |
| | I4 | | | | F | D | n | | | | | | | |
| | I5 | | | | | F | n | | | | | | | |
| | | | | | | | | | | | | | | |
| IHS | | | | | | | F | D | E | M | W | | | |

بتكون من Many Instructions

ممكن نتقل الـ Mem بس لما توصل الـ W
ما بدير نعملوا و بتحول لـ null و بنروح للـ IHS
و هيك منمنا إنه الـ instructionات اللي قبلوا خلصوا لإنه ما عليهم مشكلة
عملنالهم Flash لإنه ما بدير يكملوا

بعد ما حلينا الـ exception باستخدام register SEPC بقدر إنه يرجع للـ I3 مش
الـ add و هذا المثال خصيصا لإنه عملية الـ add ما بنزبط لإنه الناتج كبير ما بنخزن .
\* فإحنا ممكن نرجع لنفس الـ instructionات أو اللي بعدها و هنا بنقدر على الحالة .

## Slide 47 : Multiple exceptions :

حفترض إنه opcode illegal الوا ما بنقدر نعملوا في ـ E

| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | I1 | F | D | E | M | W | | | | | | | | |
| | add $X1$,$X2$,$X1$ | | F | D | Ë | M̈ | n | | | | | | | |
| | I3(bad) | | | F | Ö | Ë | n | | | F | D\* | E\* | M\* | n |
| | I4 | | | | F | D | n | | | | | | | |
| | I5 | | | | | F | n | | | دورة I3 | | | دورة I3 | |
| | | | | | | | | | | | | | | |
| IHS | | | | | | | F | D | E | M | W | | | |
| | | | | | | | | | | F | D | E | M | W |

- بالأول بنحل مشكلة الـ add قبل الـ I3 و بنرجع و بيجي دور الـ I3 بعدين

## Slide 58 : Hazards and forwarding in Daul issue static processor :

Forwoding

| Dependent load ⇃ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| add $X_{10}, X_0, X_1$ | F | D | E | M | W | | | | | | → هو حطيناهم |
| ld $X_2, 0 (X_{10})$ | | F | D | E | M | W | | | | | issue Packet بـ |
| | | | | | | | | | | | وحدة لإنه |
| | | | | | | | | | | | بعضهم على بعض يشتغلوا |
| | | | | | | | | | | | فلازم بما انه |

الـ ld بيشتغل على الـ add لازم يكون الـ ld متأخرة في issue packet 8 حجة

## Slide 58 : Load use hazard :

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ld $X_1, 0 (X_2)$ | F | D | E | M | W | | | | | | → ماينفذر نحطهم |
| sub $X_4, X_1, X_5$ | | | F | D | E | M | W | | | | بنفس الـ issue packet |
| | | | | | | | | | | | فلازم نحط الـ sub |
| | | | | | | | | | | | بـ issue packet |
| | | | | | | | | | | | ثانية بس يكون |

فيه بعد الـ ld كمان issue packet تشتني ينزبط الحل

## Slide 59 : From W in Memory pipeline to E in ALU Pipeline :

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ld $X_{31}, 0 (X_{20})$ | F | D | E | M | W | | | | | | |
| Sub $X_{31}, X_{31}, X_{21}$ | | | F | D | E | M | W | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |
| | | | | | | | | | | | |

## Slide 59 : From M in ALU pipeline to M in Memory pipeline :

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| add X31, X31, X21 | F | D | E | M | W |  |  |  |  |  |
| sd X31, 0 (X20) | F | D | E | M | W |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |  |  |  |  |

* أسرع من ←

load use hazard

لإنو الـ address

الـ sd بنحسب من

X20 وهو ما بعتمد

على الـ instruction اللي فوقها والـ X31 بنحسب بزيادة الـ cycle 3 والـ sd بدها تاخد الـ X31 وتبعتها
على الـ Memory وبتبعتها على الـ M بالـ cycle4. إذا الـ M بالـ
فبقدر أحطهم بـ issue packet بشرط نعمل forwarding.

## Slide 62 : Example of loop unrolling :

Loop:

ld   X31, 0 (X20)

add   X31, X31, X21

sd   X31, 0 (X20)

addi   X20, X20, -8

blt   X22, X20, loop

① Replicate the loop instructions n times → 4

Loop:

| | |
|---|---|
| ld X31, 0 (x20) | ld X31, 0 (X20) |
| add X31, X31, X21 | add X31, X31, X21 |
| sd X31, 0 (X20) | sd X31, 0 (X20) |
| addi X20, X20, -8 | addi X20, X20, -8 |
| x ⌊ blt X22, X20, loop | x ⌊ blt X22, X20, loop |

| | |
|---|---|
| ld X31, 0(X20) | ld X31, 0 (X20) |
| add X31, X31, X21 | add X31, X31, X21 |
| sd X31, 0 (X20) | sd X31, 0 (X20) |
| addi X20, X20, -8 | addi X20, X20, -8 |
| x⌊ blt X22, X20, loop | blt X22, X20, loop |

② Remove unneeded loop overhead

Loop:

```
ld  X31, 0 (X20)
add X31, X31, X21
sd  X31, 0 (X20)
```

```
ld  X31, 0(X20)      -16
add X31, X31, X21
sd  X31, 0(X20)      -16
```

```
ld  X31, 0 (X20)     -8
add X31, X31, X21
sd  X31, 0(X20)      -8
```

```
ld  X31, 0( X20)     -24
add X31, X31, X21
sd  X31, 0 (X20)     -24
addi X20, X20, -8    -32
blt  X22, X20, Loop
```

③ Modify instructions

Loop:

```
ld  X31, 0 (X20)
add X31, X31, X21
sd  X31, 0 (X20)
```

```
ld  X31, -16 (X20)
add X31, X31, X21
sd  X31, -16 (X20)
```

```
ld  X31, -8 (X20)
add X31, X31, X21
sd  X31, -8 (X20)
```

```
ld  X31, -24 (X20)
add X31, X31, X21
sd  X31, -24 (X20)
addi X20, X20, -32
blt  X22, X20, Loop
```

④ Rename registers

Loop:

  ld  X28, 0 (X20)                ld X30, -16 (X20)

  add  X28, X28, X21          add X30, X30, X21

  Sd  X28, 0 (X20)                Sd  X30, -16 (X20)


  ld  X29, -8 (X20)            ld X31, -24 (X20)

  add X29, X29, X21          add X31, X31, X21

  Sd  X29, -8 (X20)           Sd  X31, -24 (X20)

                            addi  X20, X20, -32

                            blt  X22, X20, Loop

★الـ Renaming رح يساعدني إني أداخل بين الـ instructions وأنا بعمل Schedule حتى أقلل من المشاكل اللي ممكن تطلوعني ( anti-dependence, output dependence, WAW, WAR )

⑤ Schedule instructions

| | ALU / branch | Load / store | cyck | ld X28, 0 (X20) |
|---|---|---|---|---|
| Loop: | addi X20, X20, -32 | ld X28, 0 (X20) | 1 | add X28, X28, X21 <br> sd X28, 0 (X20) |
| | nop | ld X29, 24 (X20) | 2 | ld X29, -8 (X20) <br> add X29, X29, X21 |
| | add X28, X28, X21 | ld X30, 16 (X20) | 3 | Sd X29, -8 (X20) |
| | add X29, X29, X21 | ld X31, 8 (X20) | 4 | ld X30, -16 (X20) <br> add X30, X30, X21 |
| | add X30, X30, X21 | Sd X28, 32 (X20) بدل 0 م | 5 | Sd X30, -16 (X20) <br> ld X31, -24 (X20) |
| | add X31, X31, X21 | Sd X29, 24 (X20) | 6 | add X31, X31, X21 |
| | nop | Sd X30, 16 (X20) | 7 | sd X31, -24 (X20) <br> addi X20, X20, -32 |
| | blt X22, X20, Loop | Sd X31, 8 (X20) | 8 | blt X22, X20, Loop |

★ نخلص تعبية الـ loads أول لأنو بفضل تنفيذهم اسرع ما يمكن لإنو الـ data الـ instructions اللي بحتاجوها الـ load بنجيبوا

★ بغير الـ offset عشان يزبط مع الـ addi :  -8 = X-32 / -16 = X-32 / -24 = X-32

# Dynamic Multiple Issue

أحدث وأقوى
ودحرفوا كهرباء أكثر

- "Superscalar" processors → dynamic
  Multiple issue
  الاسم الثاني لـ

- CPU decides whether to issue 0, 1, 2, …
  each cycle
  مثل الـ ALU
  ما بصير يعملوا
  2 instruction مع بعض
  - Avoiding structural and data hazards

- Avoids the need for compiler scheduling
  - Though it may still help
    ما بنحتاجه
    ولكن لو عمل مثلا
    Loop unrolling
    يحسن أكثر
  - Code semantics ensured by the CPU

\* حجة الـ Code مسؤولية الـ CPU هش مسؤولية الـ Compiler

# Dynamic Pipeline Scheduling

Allow the CPU to execute instructions out of order to avoid stalls

إنه إذا فيه instructions قادرين بتشوا بعض بظلم يستننوا لكن لو في بدهم instructions أخرين بقدر أخدهم فبيخدموهم

- But commit result to registers in order

حتى نحافظ على ترتيب البرنامج ببعت النتيجة

Example

الوظائفية لما register على الترتيب الأصلي للبرنامج .

```
ld    x31,20(x21)
add   x1,x31,x2
sub   x23,x23,x3
andi  x5,x23,20
```

- Can start sub while add is waiting for ld

أهم شي بالعمله

# Dynamically Scheduled CPU

التسلسل

F

I

E

W

C

بنفيدوا حسب نوعها → به

و لما يجيبهم بالترتيب
لأنه لبعرف كل Instruction
بتعتمد على شو

**In-order issue** ← **Preserves dependencies**

| Instruction fetch and decode unit |

نقعد وا فيه ال Inst لحتى يتعين
موعد تنفيذهم

| Reservation station | Reservation station | . . . | Reservation station | Reservation station |

**Hold pending operands** ←

Functional units

| Integer | Integer | . . . | Floating point | Load-store |

عدهم → بعتمد عالتصميم
مثل في Intel i7 عدهم 7

**Out-of-order execute**

CDB

**Results also sent to any waiting reservation stations** ←

بعطل زبون results الـ يكفش الأخرى الذي تنازعوا

| Commit unit |

**Reorders buffer for register writes**

**In-order commit**

**Can supply operands for issued instructions**

أي Reservation stations فيها instruction بتنتهي لأنها result الـ بتعتمد عاشي قبلها بوديها عليها

آخر مرحلة ال result بشرط
commit unit من
register file الـ

Reservation stations فيها
instruction
بتنتهي لأنها result الـ
بتعتمد عاشي قبلها بوديها عليها

CDB / commit unit      / RF      3مصادر الـ operands ←

(result) RAW وهاد بحل مشكلة

MK

# Pipeline Stages

**F**: Fetch from instr. memory (IM) to instr. queue (IQ).

**I**: Issue from IQ to reservation stations (RS), reading ready operands from register file (RF).

**E**: Execute when functional unit (FU) is free and instr. In RS has ready operands. → متغيراتها جاهزين

**W**: Write result from FU through common data bus (CDB) to reorder buffer (ROB) and RS.

**C**: Commit results in order from ROB to RF and memory → Store بال بنكتب فحا

- Loads have **FIAMWC**, stores have **FIAC**. **A**: Address calculation

E بد ال
Address نشمي
Calculation

Memory بال c بنكتب بلا

Chapter 4 — The Processor — 67

# Register Renaming

RS [OP | OP₁ | OP₂]

- Reservation stations and reorder buffer effectively provide register renaming
- On instruction issue to reservation station
  - If operand is available in register file or reorder buffer
    - Copied to reservation station
    - No longer required in the register; can be overwritten
  - If operand is not yet available
    - It will be provided to the reservation station by a function unit
    - Register update may not be required

Chapter 4 — The Processor — 68

# Examples

- Assume superscalar processor of deg
- Name dependence (WAR)

```
mul   x1,x2,x3
add   x4,x1,x5
ld    x5,16(x21)
```

بال 5 stage Pipeline
لها لما نتفذهم عالترتيب ماكان
في مشكلة بنقدم

لـ ممكن نحله بتغير اسم الـ Register
الخوف لما يصير كم ترتيب في التنفيذ وتعمل الـ ld
قبل الـ add
وهيك بصير غلط

## Output dependence (WAW)

```
mul   x1,x2,x3
add   x4,x1,x5
ld    x1,16(x21)
```

الخوف هون ثانية إنه هاي الـ instruction
اللي هي لما تسبق الـ mul
وهيك خطأ.

ممكن نحلها بإنه نغير
اسم الـ Vegister

# Speculation →

مبني على الـ dynamic Exec ..
وبنزيد عليه الـ Speculation

- "Guess" what to do with an instruction
  - Start operation as soon as possible
  - Check whether guess was right
    - If so, complete the operation
    - If not, roll-back and do the right thing
- Common to static and dynamic multiple issue
- Examples
  - Speculate on branch outcome
    - Roll back if path taken is different
  - Speculate on load
    - Roll back if location is updated

# Compiler/Hardware Speculation

- Compiler can reorder instructions
  - e.g., move load before branch
  - Can include "fix-up" instructions to recover from incorrect guess
- Hardware can look ahead for instructions to execute
  - Buffer results until it determines they are actually needed
  - Flush buffers on incorrect speculation

# Branch Speculation

جدًا مفيد لإنه التنبؤ حاليا بكون صحيح بنسبة 95% تقريبا

- Predict branch and continue issuing
  - Don't commit until branch outcome determined → ما بنعمل commit إلا بس
  
  نعمل لل Branch عشان نتأكد من التنبؤ
- **Example**: Assume a superscalar processor of degree 2 and the branch prediction is not taken. → بقدر يجيب 2 instructions Every cycle

  ```
  ld    x1,0(x20)
  beq   x1,x2,Skip
  I3
  I4
  ```

# Load Speculation

- Avoid load and cache miss delay
  - Load before completing outstanding stores
  - Predict the effective address or loaded value
  - Bypass stored values to load unit
- Don't commit load until speculation cleared
- **Example**: Superscalar of degree 3.

```
       ld    x1,0(x20) ⎤
independent ⎡ sd  x2,0(x1)   ⎥ RAW
ما بقتموا بعض ⎣ ld  x3,0(x21) → ⎦
على الأكثر
```

```
F I A M WC
F I        A C
F I
           A M WC
```

ما بتعتمدعاللي فوق .

في حالات نادرة بكونوا
لو مثلاً الـ dependent
الـ Address اللي حطيناه بـ x1 بالصدفة كان يساوي
الـ Address اللي موجود بـ x21 وهاي حالة نادرة

Speculation لو بدون
حنضطر نتأخر لنعرف
Speculation يا سبب بالـ

# Speculation and Exceptions

- What if exception occurs on a speculatively executed instruction?
  - e.g., speculative load before null-pointer check
- Static speculation
  - Can add ISA support for deferring exceptions
- Dynamic speculation
  - Can buffer exceptions until instruction completion (which may not occur)

# Exceptions Examples

- Assume superscalar processor of degree 3 with 2 address calculation units

- E1: Predict branch as not take, but resolve to taken. The ld has exception in M.

```
beq   x1,x2,L1
ld    x5,16(x21)
```

- E2: Assume first sd has exemption in C.

```
ld    x1,0(x20)
sd    x1,0(x21)
sd    x2,16(x21)
```

المفاهيم الحديثة فيها Commitunit وهي بتساعد كتير
منها إنها إنها بتحل الـ WAW والنتائج بتنخزن حسب الترتيب الأطليبتاعهم، وجود RoB بمكنا إنه نعمل
وأيضنا بمكنا نتعامل مع الـ Exception بطريقة سليمة Speculation

# Why Do Dynamic Scheduling?

- Why not just let the compiler schedule code?

- Not all stalls are predicable

  - e.g., cache misses

- Can't always schedule around branches

  - Branch outcome is dynamically determined

- Different implementations of an ISA have different latencies and hazards

# Does Multiple Issue Work?

**The BIG Picture**

- Yes, but not as much as we'd like
- Programs have real dependencies that limit ILP
- Some dependencies are hard to eliminate
  - e.g., pointer aliasing
- Some parallelism is hard to expose
  - Limited window size during instruction issue
- Memory delays and limited bandwidth
  - Hard to keep pipelines full
- Speculation can help if done well

# Power Efficiency

→ much more complex than single pipeline

- Complexity of dynamic scheduling and speculations requires power

  *more transistors more circuites*

- Multiple simpler cores may be better

ـ ما كان بحاجة نحط مروحة

و تغيرخلال ٤٠ سنة تقريبا

← مختلف عن اللي درسناه لانشذ بـ Simple Pipeline زي اللي درسناه

| Microprocessor | Year | Clock Rate | Pipeline Stages | Issue Width | Out-of-Order/ Speculation | Cores/ Chip | P |
|---|---|---|---|---|---|---|---|
| Intel 486 *Cisc* | 1989 | 25 MHz | 5 | 1 | No | 1 | 5 |
| Intel Pentium | 1993 | 66 MHz | 5 | 2 | No | 1 | 10 |
| Intel Pentium Pro | 1997 | 200 MHz | 10 | 3 | Yes | 1 | 29 |
| Intel Pentium 4 Willamette | 2001 | 2000 MHz | 22 | 3 | Yes | 1 | 75 |
| Intel Pentium 4 Prescott | 2004 | 3600 MHz | 31 | 3 | Yes | 1 | 103 |
| Intel Core | 2006 | 2930 MHz | 14 | 4 | Yes | 2 | 75 |
| Intel Core i5 Nehalem | 2010 | 3300 MHz | 14 | 4 | Yes | 2-4 | 87 |
| Intel Core i5 Ivy Bridge | 2012 | 3400 MHz | 14 | 4 | Yes | 8 | 77 |

لـ لما الـ

Pipeline بنزيد الاشي اللي Pipeline stages

بدي لعله بـ stage وحدة بصير أقل والدائرة بـ stage وحدة بتخيير أصفر وبالتالي بتقدر تشتغل على shorter clock period

Perform... يعني أعلى

وهذا كان يزيد الـ Power كثير سيطروا عليه آخر السنوات وقللوا عدد الـ Piepline stages

# Contents

مثالين مومات

modern Process ↵

استخدم من شركة Apple في بعض هواتفها و السبب، انه مصمم يصرف كهرباء أقل

# Cortex A53 and Intel i7

برضه يستخدم في high end processor الـ

PC زي الـ computer

الجيل الأول هاد

| Processor | ARM A53 | Intel Core i7 920 |
|---|---|---|
| Market بشو مستخدم | Personal Mobile Device | Server, cloud |
| Thermal design power | 100 milliWatts (1 core @ 1 GHz) | 130 Watts |
| Clock rate | 1.5 GHz | 2.66 GHz |
| Cores/Chip | 4 (configurable) | 4 → ممكن هلانلاقي 8 و 6 |
| Floating point? | Yes | Yes |
| Multiple issue? | Dynamic | Dynamic |
| Peak instructions/clock cycle | 2 → 2 instruction every cycle | 4 → 4 instruction every cycle |
| Pipeline stages | 8 ← أمنى لأنه بيشتغل s Frequency أقل | 14 |
| Pipeline schedule | Static in-order | Dynamic out-of-order with speculation |
| Branch prediction | Hybrid طرق متطورة أكثر من الـ BHT | 2-level |
| 1st level caches/core | 16-64 KiB I, 16-64 KiB D | 32 KiB I, 32 KiB D |
| 2nd level caches/core | 128-2048 KiB | 256 KiB (per core) |
| 3rd level caches (shared) | (platform dependent) | 2-8 MB |

يعني بشرحهم

§4.11 Real Stuff: The ARM Cortex-A53 and Intel Core i7 Pipelines

# ARM Cortex-A53 Pipeline

8 stages في
بأحسن الأحوال

Address generation unit

أهم اشي فيه هو الـ PC

Instruction cach الي فيه الـ Instructions الي نستخدمها

① F1 ② F2 **F** F3 F4 زيادة 2 cycles لو فيه ها لأنشاء Iss ⑥ Ex1 ⑦ Ex2 ⑧ Wr

**I → integer instructions**

**Instruction fetch & predict**

AGU + TLB
PC

Instruction Cache

Hybrid Predictor → branch لوكان

Indirect Predictor

jump لوكان

2 Instruction

D

**Integer execute and load-store**

Integer Register file

Memory Address Calculation

Issue

تتبعت الـ Inst... حسب نوعها دالـ FU مختلفة

ALU pipe 0
ALU pipe 1
MAC pipe
Divide pipe
Load pipe
Store pipe

Writeback

**Instruction Decode**

Early Decode → 13-Entry Instruction Queue → Main Decode → Late Decode

D1 ③ ④ D2 ⑤ D3

**Floating Point execute**

NEON Register file

MUL/DIV/SQRT pipe

ALU pipe

F1 F2 F3 F4 F5

**I → Floating point instructions** لوكان

لوكان Floating فـ 10 cycles or stages

Chapter 4 — The Processor—8

# ARM Cortex-A53 Performance

Handwritten annotations (Arabic/English):

بالدراسة هاي قسموا قيمة بياخذ time بال Sec و قيمه بتاخذ Instructions

∴ Secxf →cycle
∴ cycle / IC

قسموا عدد العلبي الي بنحتاجهم لتنفيذ كل الـ Instructions لهذا البرنامج يطلعوا معي هذا الـ ratio
CPI

Memory hiraray Stalls

مفروض تكون 0.5 (IPC = 2)
كل ما كانت اقل كلما كانت أفضل

المتوقع

الها سببين : يا إنه الـ Instruction في بنيها اعتماد يعطض أو إنه الـ Processor بحاول يطلب الـ data من الكاش وهيك بصير

بحتاج فترة طويلة ليتجي الـ data من الـ Mem

# Core i7 Pipeline

→ Processor much more complex
Superscalar Processor with degree 4

*ما يتغير من حاسبة لحاسبة لهذا الي بختلف من وقت*

*لوقت عدد وأحجام الـ IQ و RS و RoB*
*وحتى عدد الـ FU*

**128-Entry inst. TLB (four-way)**

**32 KB Inst. cache (four-way associative)** ← IM *ما بنعرف كم فيهم Instructions*

**16-Byte pre-decode + macro-op fusion, fetch buffer** ← F

**Instruction fetch hardware**

**18-Entry instruction queue** ← IQ *بسع 18 مجموعة → 18*16 byte*

*complex Instructions بنرجمها Simple* →

Micro-code | **Complex macro-op decoder** | **Simple macro-op decoder** | **Simple macro-op decoder** | **Simple macro-op decoder**

D → 4 instruction لـ
— decode instructions الـ
بينحطوا فيه

**28-Entry micro-op loop stream detect buffer**

Issue

**Retirement register file**

**Register alias table and allocator**

**128-Entry reorder buffer**

I ← RS

*بمرحلة الـ Issue بنجبر لكل Instruction بالـ RoB بالـ بيضمن الـ commit ترجيل النتائج الـ الأصلي حسب ترتيبو*

**36-Entry reservation station**

| ALU shift | ALU shift | Load address | Store address | Store data | ALU shift |
| SSE shuffle ALU | SSE shuffle ALU | | **Memory order buffer** | | SSE shuffle ALU |
| 128-bit FMUL FDIV | 128-bit FMUL FDIV | | | Store & load | 128-bit FMUL FDIV |

E ← تنفيذ الـ Instructions بـ (FU)
Functional unit حسب نوعهم +
بقدروا يعملوا الـ Floating Point

**512-Entry unified L2 TLB (4-way)** | **64-Entry data TLB (4-way associative)** | **32-KB dual-ported data cache (8-way associative)** | **256 KB unified l2 cache (eight-way)**

*+ أهم شي*

**8 MB all core shared and inclusive L3 cache (16-way associative)** → **Uncore arbiter (handles scheduling and clock/power state differences)**

M

*مش حنجلي نعمل بالوقت الحاضر*

Chapter 4 — The Processor — 83

# Core i7 Performance

*بتقلل من الـ Performance*

Stalls, misspeculation
Ideal CPI → = 1/4 المفروض

CPI values: 0.44, 0.59, 0.61, 0.65, 0.74, 0.77, 0.82, 1.02, 1.06, 1.23, 2.12, 2.67

(libquantum, h264ref, hmmer, perlbench, bzip2, xalancbmk, sjeng, gobmk, astar, gcc, omnetpp, mcf)

*أسوأ ما يمكن 10%*

Branch misprediction %   Wasted work %

(libquantum, h264ref, hmmer, perlbench, bzip2, xalancbmk, sjeng, gobmk, astar, gcc, omnetpp, mcf)

*٣ مرات أسرع من الـ ARM لأنه بشتغل على high frequency في dynamic exc... وكلها مع بعضها بتعطي نتائج أفضل لكن على حساب الـ Power*

Chapter 4 — The Processor — 84

Scanned with CamScanner

# Contents

# Fallacies

- Pipelining is easy (!)
  - The basic idea is easy
  - The devil is in the details
    - e.g., detecting data hazards
- Pipelining is independent of technology
  - So why haven't we always done pipelining?
  - More transistors make more advanced techniques feasible
  - Pipeline-related ISA design needs to take account of technology trends
    - e.g., predicated instructions

# Pitfalls

Poor ISA design can make pipelining harder

- e.g., complex instruction sets (VAX, IA-32)
  - Significant overhead to make pipelining work
  - IA-32 micro-op approach
- e.g., complex addressing modes
  - Register update side effects, memory indirection
- e.g., delayed branches
  - Advanced pipelines have long delay slots

# Contents

# Concluding Remarks

Pipelining improves instruction throughput using parallelism

- More instructions completed per second
- Latency for each instruction not reduced

Hazards: structural, data, control

Multiple issue and dynamic scheduling (ILP)

- Dependencies limit achievable parallelism
- Complexity leads to the power wall

## Slide 67: Dynamic Scheduling:

يعني يجيب بس instruction 1 →

**Single issue** →

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| ld X31, 20 (X21) | F | I | A | M | W* | C | | | | |
| add X1, X31, X2 | | F | I | – | – | E | W | C | | |
| Sub X23, X23, X3 | | | F | I | E | W* | – | – | C | |
| andi X5, X23, 20 | | | | F | I | – | E | W | – | C |

\* بنعمل هم I بالترتيب الأصلي تبعهم وكذلك الـ C ولكن E مش عالترتيب الجاهز ننفذ

## Slide 69:

يعني أي stage بتقدر تعمل 3 things →

\* Triple Issue : Name dependence (WAR)

نفترض إنه الـ mul تحتاج كـ 3cycles بالـ EX

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| mul X1, X2, X3 | F | I | E | E | E | W | C | | | |
| add X4, X1, X5 | F | I | – | – | – | – | E | W | C | |
| ld X5, 16 (X21) | F | I | A | M | W | – | – | – | C | |

نجمع محتويات X21 مع 16

\* الكتابة على الـ RF بتتم بالـ C مش الـ W

\* Triple Issue : output dependence (WAW)

نفترض إنه الـ mul تحتاج 3 cycles بالـ Ex

كتبت بـ ROB مش بالـ RF

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| mul X1, X2, X3 | F | I | E | E | E | W | C | | | | → هون بكتب ع X1 |
| add X4, X1, X5 | F | I | – | – | – | – | E | W | C | |
| ld X1, 16 (X21) | F | I | A | M | W | – | – | – | C | | → هون بكتب ع X1 |

وهون X بتخزن بالـ RF لإنه النتيجة X النهائية

هاي كتبت بالـ reorder buffer

RF مش بالـ

Slide 72 :

Example: Assume a superscalar processor of degree 2 and the branch Predection is not taken. ( Correct Prediction )

فحصنا C لل I3

بقسم الـcycle على

لما عملنا C اكتشفنا

1. صح انتشؤ

بمعلوم F بالـ →

2 cycle إنه توقعنا إنه

إنه not taken

الـ Branch وتوقعنا

صحيح كان.

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| ld X1 , 0 (X20) | F | I | A | M | W | C | | | | |
| beq X1 , X2, Skip | F | I | | | | | E | W | C | |
| I3 | | | F | I | E | W | | | C | |
| I4 | | | F | I | E | | W | | | C |
| ... | | | | | | | | | | |
| Skip: | | | | | | | | | | |

3 write أعمل بقدر هادي لأنه W الـ أخذنا لـ

of degree 2  بنفس الـcycle لأنه البرنامج

فعليّاً بقدر أعمل 2 من F / I / M / W / C

Functional unit بس الـ E بعنمو قديش عندي

E   A

Slid 72 :

Example: Assume a superscalar Processor of degree 2 and the Branch Predection is Not taken. (Incorrect Prediction )

بس عملنا E

اكتشفنا إنه لأ لتنبؤ

خاطئ

ما بنعملوش C إلا

بعد الـ Branch →

بنعملها C

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| ld X1 , 0 (X20) | F | I | A | M | W | C | | | | |
| beq X1 , X2, Skip | F | I | | | | | E | W | C | |
| I3 | | | F | I | E | W | | | n | |
| I4 | | | F | I | E | | W | | n | |
| --- | | | | | | | | | | |
| Skip: | | | | | | | | | F | |

## Slide 73 :

Example : Load speculation. Assume a superscalar processor of degree 3.
Predict the second load does not depend on the store.
(Correct prediction)

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| ld $X_1$ , 0($X_{20}$) | F | I | A | M | W | C | | | | |
| sd $X_2$, 0($X_1$) | F | I | | | | A | C———————————→ | | | |
| ld $X_3$, 0($X_{21}$) | F | I | A | M | W | | C | | | |

← بنكتب على Mem
← بال Commit

## Slide 73 :

Example : load speculation. Assume a superscalar processor of degree 3.
Predict the second load does not depend on the store. (Incorrect prediction)

→ اكتشفنا إنه $X_1 = X_{21}$

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| ld $X_1$, 0($X_{20}$) | F | I | A | M | W | C | | | | |
| Sd $X_2$, 0($X_1$) | F | I | | | | A | C | | | |
| ld $X_3$, 0($X_{21}$) | F | I | A | M | W | | n | F | | |

] → RAW

على فرض، انه $X_1 ≠ X_{21}$ ↰

↰ لإنها أخذنا قيمة خاطئة ⟍
→ بنغيرها مرة ثانية لإنها كانت خطأ
بالبداية وهيك بنمشي صح.

Slide 75 : Speculation and Exceptions

Example 1: Assume superscalar processor of degree 3 with 2 address calculations units.

Predict branch as not taken, but resolve to taken. The ld has exception in M.

هون بنكتشف إنوا Taken

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| beq X1,X2,L1 | F | I | E | W | C | | | | | | | |
| ld X5,16 (X21) | F | I | A | M | n | | | | | | | |

جبناها بناء على الافتراض تاكنا إنه ال beq
Not taken ←
هون اكتشفنا إنه فيه Exception

معناها لازم نروح لل EHS → Exception Handling Subroutine

بس بالتصاميم الحديثة إذا فيه Exception احنا ما بنعالجها إلا بال Commit بس لسيتشان
احنا بننفذ ال ld عالفاضي لأنه ال beq → taken فبنسبها طبيعي بنعملها null
وال Exception و كأنه ما صار.

Example 2 : Assume superscalar processor of degree 3 with 2 address calculation units.

Assume first sd has exception in C.

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| ld X1 , 0 (X20) | F | I | A | M | W | C | | | | | |
| sd X1 , 0 (X21) | F | I | A | - | - | n | | | | | |
| sd X2, 16 (X21) | F | I | - | A | - | n | | | | | |
| EHS : | | | | | | | F | | | | |

عرفنا إنه فيه exception فنشحول من C إلك null

و نروح لل EHS وبس أخلص والأمور تكون ملظمة بِرجع بعمل F لل instruction اللي عملناهم null و برجع أنفذهم مرة ثانية.

Scanned with CamScanner

# Chapter 5

# Large and Fast: Exploiting Memory Hierarchy

التركيز على الـ Memory

، الهدف اني أصمم Memory حديثة تكون كبيرة وبنفس الوقت سريعة وهذا هش موضوع سد

*Adapted by Prof. Gheith Abandah*

# Contents

حندرسهم من ناحية ←
ال Memory

صفة من صفات البرامج لها نتعامل مع الـ Memory
و بنفس الوقت موجودة بحياتنا ؟

معقد

# Principle of Locality

- Programs access a small proportion of
  their address space at any time → مثلا برنامج
  معين ما بشتغل عكل امكانياته بشتغل عشغلات معينة فيه بس.

- Temporal locality

  من كلمة time
  - Items accessed recently are likely to be
    accessed again soon → إذا اشي انت طلبته هسا في احتمال كبير انك تطلبه بعد شوي.
    - e.g., instructions in a loop, induction variables
      ↳ EX: for ( i=0 ... )↑

- Spatial locality

  جاية من كلمة space
  - Items near those accessed recently are likely
    to be accessed soon → إذا في اشي عملتله access احتمال اعمل access الـ loctations اللي قريبين منه
    - E.g., sequential instruction access, array data

Scanned with CamScanner

# Taking Advantage of Locality

- Memory hierarchy
- Store everything on <u>disk</u> or Flash
- Copy recently accessed (and nearby) items from disk to smaller DRAM memory
  - Main memory
- Copy more recently accessed (and nearby) items from DRAM to smaller SRAM memory
  - Cache memory attached to CPU

*Arabic handwritten annotations:*
أنه كلما نكون أكثر من نوع وأكثر من طبقة

إنه المعلومات تكون أكبر وبطيئة

هذا المستوى الثاني بعد الـ disk / Flash هي بنية أكثر سريعة

الشغلات اللي منها عليها والقريبة ركن مبنية سريعة

أهم ما هنا الـ disk ولما نسكر الجهاز محتوياتها ما يروح مثل البرامج اللي مخزنة بـ Smart Phones

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy —

# Memory Hierarchy

*Arabic handwritten annotations:*
لا تصميم مثل هذا التصميم بالنسبة لا Processor مادح يحصل مع بعض Fast and large

هي اللي بنستخدمها أكثر

DRAM هنا

باللخارج disk / Flash

| | volatile | | non-volatile |

- Fast →
- KB = $10^3$ — CPU { Core }
- MB = $10^6$
- GB = $10^9$
- TB = $10^{12}$
- Large → PB = $10^{15}$

CAPACITY

- eSRAM    < 1 ns
- embeded
- Cache    3 - 10 ns
- Main memory DRAM    10 - 30 ns
- Local storage disk    ~ 1 ms
- Remote secondary storage cloud storage    ~100 ms

disk بالخارج

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy —

# Memory Hierarchy Levels

*ال Memory بتتقسم ل blocks*
*blocks*
*يعني لما نجيب وحدة هنتحت لفوق*
*بنجيبوا blocks مش متفردة*

**Block** (aka <u>line</u>): unit of copying

- May be <u>multiple words</u> 64 bytes → block *ل الواحد*

If accessed data is present in
upper level → *اذا ال Processor طلب شغلات من ال Cashe بنقول عن نوع هاي ال access, اونا Hit*

- Hit: access satisfied by upper level
  - *ب ال Cashe لقينا اللي بدنا اياه*
  - Hit ratio: hits/accesses → $= \dfrac{hits}{hits + misses}$

If accessed data is absent

- Miss: block copied from lower level *الداتا اللي بدنا ياها مش موجودة بال Cashe موجودة بال lower Memory level*
  - *بناخد ال block اللي فيه ال word المطلوبة وبنسخه للـ upper Memory level*
  - Time taken: miss penalty
  - Miss ratio: misses/accesses = 1 – hit ratio
- Then accessed data supplied from upper level

*الزمن الاضافي اللي بدي منتاجه كين ما اجيب ال data*

**Processor**

**cache Mem ($)**

Data is transferred

**Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 6**

→ EX : hit ratio = 95 % → *دايما عالي*
miss ratio = 100% – 95% = 5 %

---

# Contents

# Memory Technology (2012)

- Static RAM (SRAM) → الأسرع ولكن الأغلى
  - 0.5ns – 2.5ns, $2000 – $1000 per GB  
    لل gega byte الواحدة
- Dynamic RAM (DRAM)
  - 50ns – 70ns, $10 – $20 per GB
- Flash memory
  - 5,000ns – 50,000ns, $0.75 – $1.00 per GB
- Magnetic disk
  - 5ms – 20ms, $0.05 – $0.10 per GB
- Ideal memory
  - Access time of SRAM
  - Capacity and cost/GB of disk

← Static Random Access Memory

# SRAM Technology
" SRAMs lose data content on power loss "

- Static RAM
- 6-8 transistors per bit
- Fast but not dense → كثيف
- Often has standby mode → وضعية الاستعداد

→ select line



IDT6167SA/LA
CMOS Static RAM 16K (16K x 1-Bit)

**Pin Configurations**

$2^{14} = 2^4 \times 2^{10} = 16 \text{ k bits}$

Scanned with CamScanner

1969 ←    → Dynamic Random Access Memory

# DRAM Technology →

بنعتقد،انه نخزن المعلومات
: Capacitor و بنحتاج
one transistor

## Data stored as a charge in a capacitor

- Single transistor used to access the charge
- Must periodically be refreshed
  - Read contents and write back
  - Performed on a DRAM "row"

Select ──
Storage capacitor
Data

Column
Rd/Wr
Act
Pre
Bank
Row

* بنخزن 1 transistor
1 bit بتخزن

* DRAM : Refresh Each memory cell 16 times Per Second (or more) or risk losing content.

# Classic DRAM

## Basic DRAM chip

RAS# ──→
Addr ──→
Column
Row Address Latch ─ Row → Row Address decoder
Memory array
Column Address Latch
Column ──→ Column addr decoder
CAS# ──→
Data

- DRAM access sequence
  - Put Row on addr. bus
  - Assert RAS# (Row Addr. Strobe) to latch Row
  - Put Column on addr. bus
  - Wait RAS# to CAS# delay and assert CAS# (Column Addr. Strobe) to latch Col
  - Get data on address bus after CL (CAS latency)

13

Computer Structure 2015 – System

# Classic DRAM

RAS#

CAS#

Addr ( Row 1 )—( Col.1 )          ( Row 2 )—( Col.2 )

Data          ( Data1 )          ( Data 2 )

Every access - individual

2009-2013                    ©S.Maciulevičius

- Low bandwidth

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy—

---

# Advanced DRAM Organization

- Access an entire row and save it in a **row buffer**.

- **Fast page mode**: supply successive words from the row buffer with reduced latency

RAS#

CAS#

Addr ( Row 1 )—( Col. 1 )     ( Col. 2 )          ( Col. 3 )

Data          ( Data1 )     ( Data2 )     ( Data3 )

Chapter 5 — Large and F... ...ierarchy—

Scanned with CamScanner

# Advanced DRAM Organization

**Synchronous DRAM** (SDRAM) has a counter that increments the column address using a clock signal.

متزامن

```
                              ┌──────────────┐
                              │  Row Buffer  │
                              └──────┬───────┘
                                     │
         ┌──────────────────┐  ┌─────┴──────────┐
Addr ────│ Column Address   │──│ Column Decoder │
         │ Counter/Latch    │  └────────────────┘
         └────┬────────┬────┘
              │        │
            $\overline{CAS}$   CLK
```

# Advanced DRAM Organization

- **Double data rate (DDR) SDRAM**
  - Transfer on rising and falling clock edges
- **Quad data rate (QDR) SDRAM**
  - Separate DDR inputs and outputs



SDR — Single Data Rate — 1 signal per clock cycle

DDR — Double Data Rate — 2 signals per clock cycle

QDR — Quad Data Rate — 4 signals per clock cycle

معظم الأجهزة بتستخدمها

I Port بدل ← فيه 2 Ports

فكرته، إنه في clock والداتا بتنتقل مرتين كل cycle عند الـ Raising Edge وعند الـ falling edge خلال cycle الواحدة بيتم انتقال 2 data transfers وهذا النوع بيرفع من الـ Mem الـ bandwidth تبعها كثير

# Micron 1Gb DDR-SDRAM

MT46V128M8 – 32 Meg X 8 X 4 Banks, <u>Datasheet</u>

# Micron 1Gb DDR-SDRAM

# DRAM Generations

| Year Introduced | Chip size | $ per GiB | Total access time to a new row/column | Average column access time to existing row |
|---|---|---|---|---|
| 1980 | 64 Kibibit | $1,500,000 | 250 ns | 150 ns |
| 1983 | 256 Kibibit | $500,000 | 185 ns | 100 ns |
| 1985 | 1 Mebibit | $200,000 | 135 ns | 40 ns |
| 1989 | 4 Mebibit | $50,000 | 110 ns | 40 ns |
| 1992 | 16 Mebibit | $15,000 | 90 ns | 30 ns |
| 1996 | 64 Mebibit | $10,000 | 60 ns | 12 ns |
| 1998 | 128 Mebibit | $4,000 | 60 ns | 10 ns |
| 2000 | 256 Mebibit | $1,000 | 55 ns | 7 ns |
| 2004 | 512 Mebibit | $250 | 50 ns | 5 ns |
| 2007 | 1 Gibibit | $50 | 45 ns | 1.25 ns |
| 2010 | 2 Gibibit | $30 | 40 ns | 1 ns |
| 2012 | 4 Gibibit | $1 | 35 ns | 0.8 ns |

بتقل السعر مع الزمن

بتحسن مع الزمن

# DRAM Generations

| Year | Capacity | $/GB |
|---|---|---|
| 1980 | 64Kbit | $1500000 |
| 1983 | 256Kbit | $500000 |
| 1985 | 1Mbit | $200000 |
| 1989 | 4Mbit | $50000 |
| 1992 | 16Mbit | $15000 |
| 1996 | 64Mbit | $10000 |
| 1998 | 128Mbit | $4000 |
| 2000 | 256Mbit | $1000 |
| 2004 | 512Mbit | $250 |
| 2007 | 1Gbit | $50 |



row جديد

Trac
Tcac

column جديد

# DRAM Performance Factors

- Row buffer
  - Allows several words to be read and refreshed in parallel
- Synchronous DRAM
  - Allows for consecutive accesses in bursts without needing to send each address
  - Improves bandwidth
- DRAM banking
  - Allows simultaneous access to multiple DRAMs
  - Improves bandwidth

# Increasing Memory Bandwidth

تصاميم مختلفة للـ Memory

بتحتاج وقت طبعاً

بتحتاج وقت أقل

لا بياخذ وقت بس access لا transfer سريع



a. One-word-wide memory organization

b. Wider memory organization

c. Interleaved memory organization

## To get 16-byte block:
### a. One-word wide memory
- Miss penalty = $4 \times (1 + 15 + 1) = 68$ bus cycles
- Bandwidth = 16 bytes / 68 cycles = 0.24 B/cycle

Scanned with CamScanner

# Increasing Memory Bandwidth



b. Wider memory organization

c. Interleaved memory organization

a. One-word-wide memory organization

- b. **4-word wide memory**
  - Miss penalty = 1 + 15 + 1 = 17 bus cycles
  - Bandwidth = 16 bytes / 17 cycles = 0.94 B/cycle
- c. **4-bank interleaved memory**
  - Miss penalty = 1 + 15 + 4×1 = 20 bus cycles
  - Bandwidth = 16 bytes / 20 cycles = 0.8 B/cycle

# Increasing Memory Bandwidth

- d. **DDR-SDRAM**
  - Miss penalty = 1 + 15 + 4×0.5 = 18 bus cycles
  - Bandwidth = 16 bytes / 18 cycles = 0.89 B/cycle

DIMM

# Flash Storage

Nonvolatile semiconductor storage

- 100× – 1000× faster than disk
- Smaller, lower power, more robust
- But more \$/GB (between disk and DRAM)

# Flash Types

- NOR flash: bit cell like a NOR gate
  - Random read/write access
  - Used for instruction memory in embedded systems
- NAND flash: bit cell like a NAND gate
  - Denser (bits/area), but block-at-a-time access
  - Cheaper per GB
  - Used for USB keys, media storage, …
- Flash bits wears out after 1000's of accesses
  - Not suitable for direct RAM or disk replacement
  - Wear leveling: remap data to less used blocks

# Disk Storage

Nonvolatile, rotating magnetic storage

# Disk Sectors and Access

- Each sector records
    - Sector ID
    - Data (512 bytes, 4096 bytes proposed)
    - Error correcting code (ECC)
        - Used to hide defects and recording errors
    - Synchronization fields and gaps
- Access to a sector involves
    - Queuing delay if other accesses are pending
    - Seek: move the heads
    - Rotational latency
    - Data transfer
    - Controller overhead

# Disk Access Example

Given
- 512B sector, 15,000rpm, 4ms average seek time, 100MB/s transfer rate, 0.2ms controller overhead, idle disk

Average read time
- 4ms seek time
  + ½ / (15,000/60) = 2ms rotational latency
  + 512 / 100MB/s = 0.005ms transfer time
  + 0.2ms controller delay
  = 6.2ms

If actual average seek time is 1ms
- Average read time = 3.2ms

# Disk Access Example 2

Given
- 15,000rpm, 2MB/cylinder

Sustainable peak transfer rate?

# Disk Performance Issues

- Manufacturers quote average seek time
  - Based on all possible seeks
  - Locality and OS scheduling lead to smaller actual average seek times
- Smart disk controller allocate physical sectors on disk
  - Present logical sector interface to host
  - SCSI, ATA, SATA
- Disk drives include caches
  - Prefetch sectors in anticipation of access
  - Avoid seek and rotational delay

# Contents

مبدأ مستخدم بشكل واسع بحياتنا
وبسهل علينا.

Reading

# Cache Memory

## Cache memory

- The level of the memory hierarchy closest to the CPU

المعلومات بالعربية: المهم تتكون من طبقات قريبة من الـ CPU

Given accesses $X_1, \ldots, X_{n-1}, X_n$

أبسط حالة Cache واحد



a. Before the reference to $X_n$   b. After the reference to $X_n$

How do we know if the data is present?

Where do we look?

---

# Direct Mapped Cache

هو يحدد مكان الـ block في الـ Cache

أنا أبسط أنواع الـ cache

- Location determined by <u>address</u>

- Direct mapped: only one choice

دائماً يكون لنا كل واحد فقط عنوان على حسب الـ Memory وهذا لتعقد على address

- (Block address) modulo (#Blocks in cache)



8 blocks

EX: 21 mod 8
= 5 → 101

EX: 8 = $2^3$

بنأخذ الـ block Address وينقسم على عدد الـ blocks ما يأخذها cache وبتأخذ الباقي

- #Blocks is a power of 2

- Use low-order address bits

32 blocks

least significant bits

# Tags and Valid Bits

How do we know which particular block is stored in a cache location?

- Store block address as well as the data
- Actually, only need the high-order bits
- Called the tag

*اهم* بنحتاج سي الـ high Address bit
يعني بختاج أجيب الـ bits الباقيت بس لأنه آخر 3 أنا أخذتهم أصلاً

What if there is no data in a location?

- Valid bit: 1 = present, 0 = not present
- Initially 0

لما نخزن tag و لما يطلب Processor الـ address بالـ Cache نبحص
الـ Address اللي طالبه مع الـ tag الموجود بالـ cache إذا تساووا
معناته الـ data اللي بدنا اياها موجودة بالكاش، إذا ما تساووا معناته
في هنا Miss.

لما نجيب block من
الـ Memory مش
بس بنخزن الـ data
نائته منخزن كمان
الـ Address
تائنته

# Cache Example

- 8-blocks, 1 word/block, direct mapped
- Initial state

< 32 bit > or < 1 w >

| Index | V | Tag | Data |
|---|---|---|---|
| 000 | N = 0 | | |
| 001 | N = 0 | | |
| 010 | N = 0 | | |
| 011 | N = 0 | | |
| 100 | N = 0 | | |
| 101 | N = 0 | | |
| 110 | N = 0 | | |
| 111 | N = 0 | | |

Scanned with CamScanner

# Cache Example

| Word addr | Binary addr | Hit/miss | Cache block |
|-----------|-------------|----------|-------------|
| 22 | 10 110 | Miss | 110 |

بالعشري

cache كنه مافي بالا
block أشي

| Index | V | Tag | Data |
|-------|---|-----|------|
| 000 | N | | |
| 001 | N | | |
| 010 | N | | |
| 011 | N | | |
| 100 | N | | |
| 101 | N | | |
| 110 | Y | 10 | Mem[10110] |
| 111 | N | | |

data هايك data وحطيناها بالا

وبوجدين ال Most بنحطهم بالا tag بالا cache كشنا بالمستقبل نوف شي
Significant
bit في هون

# Cache Example

| Word addr | Binary addr | Hit/miss | Cache block |
|-----------|-------------|----------|-------------|
| 26 | 11 010 | Miss | 010 |

| Index | V | Tag | Data |
|-------|-----|-----|------|
| 000 | N | | |
| 001 | N | | |
| 010 | Y=1 | 11 | Mem[11010] |
| 011 | N | | |
| 100 | N | | |
| 101 | N | | |
| 110 | Y | 10 | Mem[10110] |
| 111 | N | | |

Chapter 5

# Cache Example

| Word addr | Binary addr | Hit/miss | Cache block |
|-----------|-------------|----------|-------------|
| 22 | 10 110 | Hit | 110 |
| 26 | 11 010 | Hit | 010 |

| Index | V | Tag | Data |
|-------|---|-----|------|
| 000 | N | | |
| 001 | N | | |
| 010 | Y | 11 | Mem[11010] |
| 011 | N | | |
| 100 | N | | |
| 101 | N | | |
| 110 | Y | 10 | Mem[10110] |
| 111 | N | | |

نفس ال
طول اللي
بدنا ياها

# Cache Example

| Word addr | Binary addr | Hit/miss | Cache block |
|-----------|-------------|----------|-------------|
| 16 | 10 000 | Miss | 000 |
| 3 | 00 011 | Miss | 011 |
| 16 | 10 000 | Hit | 000 |

| Index | V | Tag | Data |
|-------|---|-----|------|
| 000 | Y | 10 | Mem[10000] |
| 001 | N | | |
| 010 | Y | 11 | Mem[11010] |
| 011 | Y | 00 | Mem[00011] |
| 100 | N | | |
| 101 | N | | |
| 110 | Y | 10 | Mem[10110] |
| 111 | N | | |

ال وهو اللي فيه
نفس اللي أنا طالبه

# Cache Example

| Word addr | Binary addr | Hit/miss | Cache block |
|-----------|-------------|----------|-------------|
| 18        | 10 010      | Miss     | 010         |

لإنطاك وهو اللي موجود فيه مش نفس اللي أنا طالبه

| Index | V | Tag | Data |
|-------|---|-----|------|
| 000 | Y | 10 | Mem[10000] |
| 001 | N |    |            |
| 010 | Y | 10 | Mem[10010] |
| 011 | Y | 00 | Mem[00011] |
| 100 | N |    |            |
| 101 | N |    |            |
| 110 | Y | 10 | Mem[10110] |
| 111 | N |    |            |

# Address Subdivision

لو بدنا نختار منهم بنحتاج  
2 bits  
00  
01  
10  
11

تصميم الـ cache منطقي أكثر

4 bytes



بيبتدي من الـ CPU ← Address (showing bit positions) <64> bits

63 62 .... 13 12 11 ....2 1 0

within word ماينستخدمش ← Byte offset

Hit — Tag — 52 — 10

2¹⁰ = 1024 هي الـ بينستخدم لنختار وحدة ما من الـ 1024 location

Data → بتروح للـ CPU

Index

Cache blocks = 1024 (Index 0, 1, 2 ... 1021, 1022, 1023)

Valid — Tag — Data

هذا في الواقع ← Memory مرتبة بطريقة انه نحط جنبها data و valid و tag وغالباً تصنع من SRAM حتى تكون سريعة ولنكون قريبة من الـ processor

52 — 32 bit = 4 byte = 1 w

اللي ظاهر من الـ 64

# Example: Larger Block Size

أكبر من lw

- 64 blocks, 16 bytes/block
  - To what block number does address 1200 map?
  
  هذا الـ address كامل 1200 فبدي الاقي الـ block فبقسم على 16 Address
  
  Block address = $\lfloor 1200/16 \rfloor$ = 75
  
  Block number = 75 modulo 64 = 11 → Location 11 بالـ cache

| Tag | Index | Offset |
|---|---|---|
| 22 bits | 6 bits | 4 bits |

63          10  9          4  3          0

54

< 64 bit >

$2^6 = 64$

بمكونا نختار 1 من 16

$2^4 \cdot 16 = \lg 16 = \lg_2 2^4 = 4$

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 42

# Block Size Considerations

لأنه لما يصير عندك Miss الـ lw

- Larger blocks should reduce miss rate → بنجيب block كامل فالبتالي
  - Due to spatial locality spatial locality بستفيد من الـ near by words وإنه الـ Processor عالغالب حيطلب الـ
  اللي قبلها أو بعدها، فلما أجيب الـ block كامل بغطي عحالي Some misses بالمستقبل القريب
- But in a fixed-sized cache
  - Larger blocks ⇒ fewer of them
    - More competition ⇒ increased miss rate ← كل ما أكبر الـ block
    الواحد بصير عندي عدد blocks أقل
  - Larger blocks ⇒ pollution يعني بنزيد عالـ data بـ Cache ما بنحتاجها
- Larger miss penalty
  - Can override benefit of reduced miss rate
  - Early restart and critical-word-first can help

← بمجرد ما توطّل الـ Critical word بجيبها/forwarding الـ
Processor ← بدل ما تنتهي الـ فما تنشئ الـ block كله data
critical word حسب ترتيبوا بالمثال بنعت الـ block block كله
جابوا أول اشي بعدين بنطلب باقي الـ block

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 43

Scanned with CamScanner

# Block Size Considerations

10%

Miss rate

5%

0%

16    32    64    128    256

Block size
Bytes

4K ↓ حجم cache

16K

64K

256K

كل ماكان أقل أحسن ← كل ماكان

لـ بنزل بالبداية عشان لما يكبر الـ block size بنستفيد
أكثر من الـ spatial locality لعدم هيك بصير يضرنا الـ
block cache فبرجع يزداد والعدد الأقل من الـ Pollution
الـ Miss rate

# Cache Misses

- On cache hit, CPU proceeds normally
- On cache miss → لما يكون في cache miss الـ data ما بتكون موجودة بالـ cache
  - Stall the CPU pipeline وبدنا زمن أطول لنجيبها فنعملوا
  - Fetch block from next level of hierarchy لـ بطلب من اللي تحتها العيت بطلب من اللي تحت تجنا هيك
  - Instruction cache miss
    - Restart instruction fetch بطلب الـ instruction مرة ثانية
  - Data cache miss
    - Complete data access

Instruction store / load بتكون دلا لو data miss

Scanned with CamScanner

# Writing to the Cache

cpu_req.addr
**(showing bit positions)**

31 ... 14 13 ... 4 3 2 1 0

cpu_req.data

الـ address اللي هيطلع من الـ cpu

كشان خبه
1024
Data entities
$2^{10} = 1024$

الباقي
$32 - 10 - 2 - 2 = 18$

لما نكتب
يستخدموا عال cache
بنجيبها ا

18
Tag

10
Index

2 Byte offset

18 bits

V D Tag

1024

18

=

Hit

Block offset

يبقسموا
4w

لـ الهدف منهم يختار من
الـ 4w

mem_data.data 128 bit

Mux    Mux    Mux    Mux

Data Write

128

* (in grid)

Data
1024
entries

128 bit = 16 byte = 4w

Data Read

Mux

32

Data

انا احتاج الـ Prosessor يختار 1w
byte منها w

to write to the cache

بنستخدمه لما يصير في cache Miss فينجيب الـ data من الـ Memory

كشان لما بدنا نكتب حنكتب عجزة من الـ block مش كله فبداية بقرأ الـ block كامل و ببحثله عال Mux واللي ما نغيرو بدخلوا زي ما هم والـ Mux المسؤول عن الـ word اللي بدها تتغير فـ سيوخاذهم من اللي مخزن بالـ cache و الـ 1w اللي حيتم الكتابة عليها بتيجي من الـ Prossesor

**Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 46**

* الكتابة بالـ cache إلها 3 مصادر 1- من الـ Processor
2- من الـ Memory
3- من الـ cache للاشي اللي مابه يتغير

# Write-Through

( في أكثر من طريقة للتعامل مع الـ write )

## Slide 43 : Block Size Considerations

< B >

< 2B >

الـ block size → 2     اذا ضاعفت الـ block Size →

صغير فبكون عدد الـ blocks    وبنفس حجم الـ Memory

كبر         اللي عندك بصير عدد الـ blocks

( هذه أفضل )        أقل

miss Penalty →

↑   latency   ↑↑↑↑     ↑   latency     ↑↑↑↑↑↑↑↑

Req      Data Tx     Req      Data Tx

رح يزيد الزمن

Miss

Penalty

→ latency ( ثابت )

B Size

# Write-Through

هون نقلك يكتب عال cache
سس . بس ما بياخذ بعين
لاعتبار إنه ممكن
نفقد ال data
وتصير معلومات
ال $ مختلفة
عن معلومات
ال M

هم الprocessor بده يكتب شيء وال address اللي بده يكتب كليه

**On data-write hit**, could just update the block in cache

موجود بال cache

■ But then cache and memory would be inconsistent

بتصير المعلومات بال cache مختلفة عن اللي بال Memory

و أول ما أحل

Write through: also update memory →

كشان يحلوا المشكلة اللي فوق.

مكان موقت لحوظ ال data بال

But makes writes take longer

■ e.g., if base CPI = 1, 10% of instructions are stores, write to memory takes 100 cycles

Effective CPI = $1 + 0.1 \times 100 = 11$

Solution: write buffer → عبارة عن Memory سريعة

■ Holds data waiting to be written to memory

■ CPU continues immediately

بظله يكتب عال cache
و الـ write buffer لحد

Only stalls on write if write buffer is already full →

ما يتعطل إلا ما لحق الـ write buffer على الـ cpu ونقبت خيلما يرحلوهم لـ Memory

# Write-Back

write
through → الـ

Alternative: On data-write hit, just update
the block in cache → فقط بنكتب عالـ block
اللي بال Cache

O=1

- Keep track of whether each block is dirty

منيح بس
المساوئ,

When a dirty block is replaced

- Write it back to memory

- Can use a write buffer to allow replacing bloc
  كبشان لحل ← to be read first
  مشكلة البطء,

# Write Allocation

الـ data مش موجودة بالـ Cache

What should happen on a write miss?

Alternatives for write-through

①- Allocate on miss: fetch the block

②- Write around: don't fetch the block (No Allocate)

  - Since programs often write a whole block before
    reading it (e.g., initialization)

For write-back

- Usually fetch the block

write
Allocate
لاصقها بسرعة

# Example: Intrinsity FastMATH

مثال: الـ micro processor اللي بالتلاجة او لابيتة او هيك

## Embedded MIPS processor

- 12-stage pipeline في الحالة 2 caches للـ I و D فبنفس الـ cycle بيعمل load/store و F
- Instruction and data access on each cycle م

Split cache: separate I-cache and D-cache

الاسم الأصح
- Each 16KB: 256 blocks × 16 words/block اسم ثاني
- D-cache: write-through or write-back → بختاروه حسب التصميم

SPEC2000 miss rates

Instruction cache
- I-cache: 0.4% اكثر hits كليه اللي بتتعامل فال data موجودا بالعادة اكثر من الـ instruction اللي بنخذها الـ Processor
- D-cache: 11.4%
Data cache
- Weighted average: 3.2%

* الـ I-cache الـ Processor بس بقرأ منه فهو لا write-through و لا write-back لانه ماني write بالقوة بس read

لو بخذ بين الاختبار النظار

# Example: Intrinsity FastMATH



32 bit ← Processor اللي بيجي من

**Address (showing bit positions)**

31 ··· 14 13 ··· 6 5 ··· 2 1 0

32-8-4-2=18

Hit

Tag 18 | 8 | 4 Byte offset | Data

Index

16 word لنختار واحد من

256 بنستخدمه لنختار واحد من Block offset

18 bits | 16 word = 64B = 512 bits

V Tag | Data

256 entries

18 | 32 | 32 | 32

= 

Mux ← لنختار 1 من

16 word لنختار 1 من

4 bits بنحتاج

لنختار واحد من 256 بنحتاج ③ 2 = 256

Scanned with CamScanner

## Slide 47 : Write through :

P → لما ال Processor يكتب بيكتب على ال cache ← وبكتب عال Memory فبتحصل فيها inconsistent

$ → سريع

M → بطيئة

## * Write buffer :

P

①

write buffer → سريع    $ hit

②

M

* ما بكتب عال M مباشرة ، أنا بدي أوصل ال data ال Memory وبدش ياها تضيع على فال data اللي بدي ياها تروح ال Memory ببعت نسخة من اللي حكتبها بال Cache ال Write buffer وهي عبارة عن Memory سريعة فبعد هيك ال Processor بكمل شغله و ال write buffer على راحته بكتب بال Memory وهاد بحلنا مشكلة ال Performence

* ال write buffer لازم يكون كبير كفاية إنا إذا ال Processor بعت Many store instructions ودا بعني يتخزنوا بال write buffer ومع الوقت يترحلوا ال Memory .

* بترهق ال Memory لأنه ال Processor بشتغل كتير وبنفذ instructions Many وكل ال store instructions باللي بده ينفذهم رح يوخذوا ال Memory فبترهق ال Memory ، accesses، فبنروح ل Write-Back

## Slide 48 : write-Back :

on hit → Processor write just on cache ( سريع )

P

①

$

②    ③

write back   M

* لما يصير في Miss وبه يعمل replacement : يعني بده الاشي اللي يجيبه من ال Mem لتبدل ال block موجود بال cache وفيه data فبتشان ما تضيع ال Data قبلما أعمل replacement أنقل ال block اللي بدي بال cache وعليه data جديدة replacement ال Memory

* بال Miss بتكون بطيئة لأنها بتمر بمرحلتين وال M بطيئة جزأ.

## Slide 48 :

* لأجل مشكلة البطء في الـ Write-Back :

① لو كان hit بضل يقرأ ويكتب من الـ cache بسرعة وما يروح الـ M

② إن صار Miss كخطوة أولى بنقل الـ block اللي بدنا نعمله replacement الـ write buffer, إلى الـ cache وهاي عملية سريعة بعدين بنعمل replacement بعد ما بلاش الـ P خلص بياخد اللي بدو ياه . صح بطيئة لكن هي مش مشكلة بعدين بعلان راحتنا بنعمل write back

[Diagram: P box → Miss, Write Buf, $ , M, slow مو مشكلة, بتيانق P مش مشكلة]

## Slide 49 : write Allocate:

read- miss. Allocate

Write ──┬── Hit ──→ write through
        │        └→ write back
        │
        └── Miss ──→ write Allocate
                 └→ write No allocate

### Allocate

[Diagram: P box, 1)miss, 3)write, $, 2)Allocate, M]

### No Allocate

[Diagram: P box, 1)miss, $, 2)write, M]

لما يصير في Miss الـ data مو موجودة بالـ M والمقر النواتي نحوا هو الـ Memory فبكتب مباشرة بالـ Memory

لما يصير في Miss بتجيب الـ data من الـ Mem بعدين بتكتب عليوا

* الأفضل الـ Allocate لإنه بتخليني أستفيد من الـ spatial locality

* المرحلة الأولى بنفضل الـ No Allocate لإنه أسرع بالبداية . بس الـ Allocate عنه بعد نظر

# Contents

---

$$CPI_{total} = CPI_{base} + \frac{Stalls}{Instructions}$$

# Measuring Cache Performance

بتكون من جزئين →

- Components of CPU time

    ① Program execution cycles   الزمن اللي بنقضيه الـInstructions داخل الـ pipeline (F I E W C)

سلوكيات   miss   hit أو   F I A M W C   ← Includes cache hit time

    ② Memory stall cycles

       Mainly from cache misses → بسبب الـ Missis

- With simplifying assumptions:

كل ما كان أعلى كلها الـ Stall أعلى →

$$\text{Memory stall cycles}$$

حسب المعطيات بنختار أي معادلة استخدم

كل ما كان أعلى الـ Stall أعلى ←     $\frac{Missis}{hits+Missis}$     كل ما كانت أعلى كل ما الـStall أعلى

$$= \frac{\text{Memory accesses}}{\text{Program}} \times \text{Miss rate} \times \text{Miss penalty} \quad ①$$

$$= \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Misses}}{\text{Instruction}} \times \text{Miss penalty} \quad ②$$

هو هذا بين كلا اللي نختزناه قديم في Missis

Scanned with CamScanner

# Cache Performance Example

- Given
  - I-cache miss rate = 2%
  - D-cache miss rate = 4%
  - Miss penalty = 100 cycles → اذا كان في Miss بتخذ 100 cycles
  - Base CPI (ideal cache) = 2
  - Load & stores are 36% of instructions
- Miss cycles per instruction
  - I-cache: $0.02 \times 100 = 2$
  - D-cache: $0.36 \times 0.04 \times 100 = 1.44$
- Actual CPI = $2 + 2 + 1.44 = 5.44$
  - Ideal CPU is $5.44/2 = 2.72$ times faster

# Average Access Time

→ Fast
- Hit time is also important for performance
- Average memory access time (AMAT)
  - AMAT = Hit time + Miss rate × Miss penalty
    ↓
    Slow
- Example          → $f = \frac{1}{T} = 1$ GHz       miss time = hit time + miss penalty
  - CPU with 1ns clock, hit time = 1 cycle, miss          = $1 + 20 = 21$
    penalty = 20 cycles, I-cache miss rate = 5%
  - AMAT = $1 + 0.05 \times 20 = 2ns$
    - 2 cycles per instruction

# Performance Summary

When CPU performance increased
- Miss penalty becomes more significant

Decreasing base CPI
- Greater proportion of time spent on memory stalls

Increasing clock rate ← الـ Processor أسرع
- Memory stalls account for more CPU cycles

Can't neglect cache behavior when evaluating system performance

---

# Associative Caches

① Fully associative
- Allow a given block to go in any cache entry
- Requires all entries to be searched at once
- Comparator per entry (expensive) →

② *n*-way set associative
- Each set contains *n* entries
- Block number determines which set
  (Block number) modulo (#Sets in cache)
- Search all entries in a given set at once
- *n* comparators (less expensive) →

# Associative Cache Example

**① Direct mapped**

Block # 0 1 2 3 4 5 6 7

Data

Tag — 1, 2

Search → 1 comparision

12 mod 8 = 4

**2-way**
**② Set associative**

4 sets, 8 blocks
Set # 0 1 2 3

Data

Tag — 1, 2

Search → Search for two blocks

12 mod 4 = 0

**③ Fully associative**

Data

Tag — 1, 2

Search → ↑↑↑↑↑↑↑↑

ممكن يكون بأي وكرة خلازم نعمل Search لكل الـ tags

block address Number 12 ← والـ والد 3 أنواع من وين نقدر علىه يقدر
الـ Caches.

Associativity → يعني فيه ممكن نزيد الـ

# Spectrum of Associativity

- For a cache with **8 entries** = 8 blocks

**One-way set associative**
**(direct mapped)**

Block | Tag | Data
0
1
2
3
4
5
6
7

**Two-way set associative**

Set | Tag Data | Tag Data
0
1
2
3

8/2 = 4 sets

**Four-way set associative**

Set | Tag Data Tag Data Tag Data Tag Data
0
1

→ 8/4 = 2 sets

**Eight-way set associative (fully associative)**

Tag Data Tag Data Tag Data Tag Data Tag Data Tag Data Tag Data Tag Data

8/8 = 1 sets = Fully Associative

# Associativity Example

- Compare 4-block caches
  - Direct mapped, 2-way set associative, fully associative
  - Block access sequence: 0, 8, 0, 6, 8

*(handwritten: طلعنم جاهزات ← cpu ل)*

## Direct mapped

| Block address | Cache index | Hit/miss | Cache content after access 0 | 1 | 2 | 3 |
|---|---|---|---|---|---|---|
| 0 | 0 | miss | Mem[0] | | | |
| 8 | 0 | miss | Mem[8] | | | |
| 0 | 0 | miss | Mem[0] | | | |
| 6 | 2 | miss | Mem[0] | | Mem[6] | |
| 8 | 0 | miss | Mem[8] | | Mem[6] | |

*(handwritten: بطلعه من : Block Address mod 4)*

*(handwritten: 5 accesses وفضه سبب لانه في)*

*(handwritten: 5 missis)*

---

*(handwritten: لتحسين Performance الـ)*

# Associativity Example

## 2-way set associative

| Block address | Cache index | Hit/miss | Cache content after access Set 0 | | Set 1 | |
|---|---|---|---|---|---|---|
| 0 | 0 | miss | Mem[0] | | | |
| 8 | 0 | miss | Mem[0] | Mem[8] | | |
| 0 | 0 | hit | Mem[0] | Mem[8] | | |
| 6 | 0 | miss | Mem[0] | Mem[6] | | |
| 8 | 0 | miss | Mem[8] | Mem[6] | | |

*(handwritten: Block Address mod 2    بطلعه من)*

*(handwritten: 4 missis)*

## Fully associative

| Block address | | Hit/miss | one set Cache content after access | | | |
|---|---|---|---|---|---|---|
| 0 | | miss | Mem[0] | | | |
| 8 | | miss | Mem[0] | Mem[8] | | |
| 0 | | hit | Mem[0] | Mem[8] | | |
| 6 | | miss | Mem[0] | Mem[8] | Mem[6] | |
| 8 | | hit | Mem[0] | Mem[8] | Mem[6] | |

*(handwritten: كل الـ cache مفتوح إلي فبقى استخدم أي index)*

*(handwritten: 3 missis)*

# How Much Associativity

Increased associativity decreases miss rate → Performance بالتالي التحسن بيحيرلك

- But with diminishing returns

- Simulation of a system with 64KB → الـ cache حجم

D-cache, 16-word blocks, SPEC2000 → benchmark

↳ 16×4 = 64 bytes

cache blocks = $\frac{64\ KB}{64}$ = 1K = 1024 block

Direct Mapped cache
- 1-way: 10.3%  ⎤
- 2-way: 8.6%  ⎟ 1.7 %
- 4-way: 8.3%  ⎟ 0.3 %
- 8-way: 8.1%  ⎦ 0.2 %

لأ أحسن واحد

Miss Rate

* مش دائماً بنستفيد لما نروح لـ higher Associativity

1S = 4 blocks
→ 256 Sets

/ ( 4-way ) → 4 comparators كأنه في

# Set Associative Cache Organization



Address
31 30 ··· 12 11 10 9 8 ··· 3 2 1 0

22          8  ↘ $2^8$ = 256

Tag
Index

* لو صار Miss فلازم تجيب الـ data من الـ Memory وتحطها بالـ block

| Index | V Tag Data | V Tag Data | V Tag Data | V Tag Data |
| 0 | | | | |
| 1 | | | | |
| 2 | | | | |
| 253 | | | | |
| 254 | | | | |
| 255 | | | | 22   32 |

256 × 4
= 1024 block
موزعين على 256 set

مستحيل يجينا 1 على الأربعة بس بكون 1 على واحد منهم أو كلهم أصفار

4-to-1 multiplexor

Hit     Data

block بنطلع أي hit هذا على أي

# Replacement Policy

- Direct mapped: no choice
- Set associative
  - Prefer non-valid entry, if there is one
  - Otherwise, choose among entries in the set
- Least-recently used (LRU) ← Replacement الشي الي من زمان ما استخدمناه و الي بنغيره
  - Choose the one unused for the longest time
  - hardware → Simple for 2-way, manageable for 4-way, too hard beyond that أصعب hardware لا أو أكثر بصير ال
- Random
  - عشوائي Gives approximately the same performance as LRU for high associativity

# Multilevel Caches

- Primary cache attached to CPU = L1
  - Small, but fast
- Level-2 cache services misses from primary cache = L2 = Secondary
  - Larger, slower, but still faster than main memory
- Main memory services L-2 cache misses
- Some high-end systems include L-3 cache
  
  Secondary ل .غرفة السامح

Scanned with CamScanner

# Multilevel Cache Example

- Given
  - CPU base CPI = 1, clock rate = 4GHz $\rightarrow \frac{1}{f} = 0.25$ ns
  - Miss rate/instruction = 2%
  - Main memory access time = 100ns = 400 cycles
- With just primary cache
  - Miss penalty = 100ns/0.25ns = 400 cycles
  - Effective CPI = 1 + 0.02 × 400 = 9

# Example (cont.)

- Now add L-2 cache
  - Access time = 5ns
  - Global miss rate to main memory = 0.5%
- Primary miss with L-2 hit
  - Penalty = 5ns/0.25ns = 20 cycles
- Primary miss with L-2 miss
  - Extra penalty = 500 cycles $\rightarrow$ 400 بدل
- CPI = 1 + 0.02 × 20 + 0.005 × 500 = 3.9
- Performance ratio = 9/3.9 = 2.3

1 level of cache     2 level of cache

# Multilevel Cache Considerations

Primary cache

- Focus on minimal hit time

L-2 cache

- Focus on low miss rate to avoid main memory access
- Hit time has less overall impact

Results

- L-1 cache usually smaller than a single cache
- L-1 block size smaller than L-2 block size

## Slide 55 : AMAT:

Find the average Memory Access time of a Memory hierarchy with the following specifications.

| Memory level | Hit Time | Miss Rate |
|---|---|---|
| L1 cache | 1 cycle | 6.0% |
| L2 Cache | 10 cycles | 4.0% |
| L3 Cache | 20 cycles | 2.0% |
| Main Memory | 250 cycles | 0% |

P → L1 → L2 → L3 → M

AMAT = Hit time + Miss rate × Miss Penalty

مو معطاه بشكل صحيح بس ال Main Memory ↳
ماعليها ال Miss فا average = 250 فلو
بما نروح من ال L3 ال Memory رح نحتاج 250
فال Miss penalty ل L3 = 250 ، فحل السؤال حيكون من الأسفل الحت الأعلى.

$AMAT_{L3} = 20 + 0.02 \times 250 = 25$ cycles
$AMAT_{L2} = 10 + 0.04 \times 25 = 11$ cycles
$AMAT_{L1} = 1 + 0.06 \times 11 = 1.66$ cycles
$\therefore AMAT = 1.66 \rightarrow$ ال يشوفه ال processor

## Slide 64 : LRU Replacement policy :

4-way associative Cache, initially empty

| Address | 1 | 2 | 3 | 4 | 1 | 5 | 3 | 6 | | Information |
|---|---|---|---|---|---|---|---|---|---|---|
| MRU | 1 | 2 | 3 | 4 | ① | 5 | ③ | 6 | ← | maintained for LRU |
| | | | 1 | 2 | 3 | 4 | 1 | 5 | 3 | replacement for each |
| | | | | 1 | 2 | 3 | 4 | 1 | 5 | Set |
| LRU | | | | | 1 | 2 | 3 | 4 | 1 | |

→ hit    → hit

يعني بحتاج 4 register احتفظ بيهم مين
ال MRU مين ال LRU

\* MRU = Most recently Used

\* LRU = Least Recently Used

# Interactions with Advanced CPUs

*dynamic execution*

Out-of-order CPUs can execute instructions during cache miss

- Pending store stays in load/store unit

- Dependent instructions wait in reservation stations

    Independent instructions continue

Effect of miss depends on program data flow

- Much harder to analyse

- Use system simulation

# Interactions with Software

محتلفة 2 Algorthims

**Misses depend on memory access patterns**

- Algorithm behavior
- Compiler optimization for memory access



$\leftarrow \dfrac{o(k)}{k}$

$\dfrac{o(k \lg k)}{k}$

برنامج الداتا

$4k \rightarrow \lg_2 4K = 12$

For large data
الـ Radix بتبين أفضل
ولكن لما نحسب
الـ clockcycle
بطلع الـ Quick
Sort أحسن
وفعلاً هو الأحسن
رغم إنه يحتاج
عدد Instructions
أكثر

فيستخدم Cache missis الى كبير الـ Radix الـ مع
large Memory عشان هيك بطيء

سيستخدم حجم أقل من الـ Mem والـ cache missis أقل

الـ Radix بدو large Memory وليبي فيها الـ data اللي بتيجي بعدين لما يجي بده يقراها بقراها بشكل متعب من مكان لمكان لأنه الـ Cache ما بوسع كل هالداتا بحيرن عنده alot of cache missis بينما الـ Quick فيه شغل على منطقة منطقة at a time ففيه عنده betteMemory access

# Contents

يعني اللي بنحطه بالـ Mem
نود نقدر نقراه

# Dependability

نذكر هل النظام شغال ولا خربان ←



Service accomplishment
Service delivered
as specified

يعني بيشغذ المطلوب ←
بشكل صحيح

✓
ok

Restoration ← لما نطلع
العطل اللي
كان

Failure

Fault: failure of a
component {→ hardware
→ Software

■ May or may not lead
to system failure

يعني هذا النظام اللي بيخدمنا يصير فيه ←
كطلامين سواء بال Software
أو ال hard ware

X

Service interruption
Deviation from
specified service

في مشكلة ←

ممكن لو صار fault مايقذي ل Failure ←
مثل لو صار تعطل بالكيبورد فهو مادخله أو لو تعطل
أو لو كان في redundancy في بديل لأشي لو عطل فما بصير ال System
ما بنستخدمه harddisck
Failure

# Dependability Measures

معتمد عليها ←

Reliability: mean time to failure (MTTF)

Service interruption: mean time to repair (MTTR)

Mean time between failures

■ MTBF = MTTF + MTTR

Availability = MTTF / (MTTF + MTTR)

نحسن availability
تكون واحد ٠ ١٠٠% ←
بس أحسن نظام
بكون اله
٩٩.٩٩٩...%

متاح ←

Improving Availability

■ Increase MTTF: fault avoidance, fault tolerance, fault
forecasting

■ Reduce MTTR: improved tools and processes for
diagnosis and repair

# The Hamming SEC Code → تصحيح الأخطاء

## Hamming distance → بقدر أمثل الأرقام

- Number of bits that are different between two bit patterns

  نقل
  ← بتقل معينة
  $\begin{matrix} 1. & 0 & 0 & 0 & 1 \end{matrix}$ → different = 1 bit
  $\begin{matrix} 2. & 0 & 1 & 0 \end{matrix}$ → " = 2 bit
  1 = Minimum السؤال

- Minimum distance = 2 provides single bit error detection → بقدر أمثل الأرقام واكتشف الخطأ
  - E.g. parity code

- Minimum distance = 3 provides single error correction, 2 bit error detection

بقدر أمثل الأرقام والكتشف الخطأ وأحله

# Encoding SEC

- To calculate Hamming code:
  - Number bits from 1 on the left
  - All bit positions that are a power 2 are parity bits
  - Each parity bit checks certain data bits:

| Bit position | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Encoded data bits | p1 | p2 | d1 | p4 | d2 | d3 | d4 | p8 | d5 | d6 | d7 | d8 |
| p1 | X | | X | | X | | X | | X | | X | |
| p2 | | X | X | | | X | X | | | X | X | |
| p4 | | | | X | X | X | X | | | | | X |
| p8 | | | | | | | | X | X | X | X | X |

odd → الا
Parity bit coverate

# Decoding SEC → Single Error Correction

Value of parity bits indicates which bits are in error

- Use numbering from encoding procedure
- E.g.

    Parity bits = 0000 indicates no error

    Parity bits = 1010 indicates bit 10 was flipped

## Example:

- Parity What will be stored for 1001 1010? (9, 10)
- If you read 0111 0010 1110, is there error? Correct it.

MK MK

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy —

Single Error Correction/
Double Error detection

# SEC/DED Code

- Add an additional parity bit for the whole word ($p_n$) → اضافة إلى الك Parity Bit

- Make Hamming distance = 4
- Decoding: ويقدر أصلحهم بس في مرات يكونوا Single Errors الـ لغير, انضخالبيت

    $P_1 P_2 P_4 P_8$

    - Let H = SEC parity bits أصلحها بقدر ما بس يكونوا double error

        H = 0, $p_n$ even, no error في error

        H ≠ 0, $p_n$ odd, correctable single bit error فبيعطي رسالة للمستخدم إنه في error لقدر المنشوم بس ما بقس أصلحه

        H = 0, $p_n$ odd, error in $p_n$ bit

        H ≠ 0, $p_n$ even, double error occurred

    Correctable بس 1 error ← لأنه H بأش كمكانه

    هدا الـ error لـ دال P ما بقوما data الـ تكون صحيحة

ECC DRAM uses SEC/DED with 8 bits protecting each 64 bits

ما بنقدر نصلح الـ double بس بنكتشفه

8 Parity Bit
H = 7
$p_n$ = 1

└ Error correction code

Chapter 5 — Large and Fast...

Scanned with CamScanner

# Contents

# RAID

رخيص

Redundant Array of Inexpensive (Independent) Disks

- Use multiple smaller disks (c.f. one large disk)
- Parallelism improves performance   م تكون مكونة data كشأن الا
  بأكثر من
- Plus extra disk(s) for redundant data storage   disck

- Provides fault tolerant storage system →   يشتغل الأخطاء
  وما بتأثر جيوا

- Especially if failed disks can be "hot swapped"

① RAID 0

وهو الـ Server شغال بتقدر إنك تشيل
الـ component المعطلة وتحط واحد
جديد

- No redundancy ("AID"?)
  - Just stripe data over multiple disks
- But it does improve performance

# RAID 1 & 2

## RAID 1: Mirroring
- N + N disks, replicate data
  - Write data to both data disk and mirror disk
  - On disk failure, read from mirror

## RAID 2: Error correcting code (ECC)
- N + E disks (e.g., 10 + 4)
- Split data at bit level across N disks
- Generate E-bit ECC $\rightarrow$ P, P2 P4 P8
- Too complex, not used in practice

# RAID 3: Bit-Interleaved Parity

- N + 1 disks
  - Data striped across N disks at byte level
  - Redundant disk stores parity
  - Read access  الداتا موزعة فنازع
    - Read all disks  نقرا من كل الـ disks
  - Write access
    - Generate new parity and update all disks
  - On failure  لو تهطل واحد بجيبه من الـ Parity
    - Use parity to reconstruct missing data  بعدين
- Not widely used

# RAID 4: Block-Interleaved Parity

### N + 1 disks

- 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
- 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

* الأوفر RAID 3 / RAID4 من الـ RAID2 لإنهم أرخص

MK

---

# RAID 3 vs RAID 4 → مستخدم أكثر من RAID 3



N= 4

New Data    1. Read 2.Read 3.Read

| D0' | | D0 | D1 | D2 | D3 | P |

+ XOR

| D0' | D1 | D2 | D3 | P' |

4.Write                    5.Write

N-1 reads
2 writes

N= 4

New Data 1.Read                    2.Read

| D0' | | D0 | D1 | D2 | D3 | P |

+ XOR

+ XOR

| D0' | D1 | D2 | D3 | P' |

3.Write                    4.Write

2 reads
2 writes

MK

# Slide 68 : Multilevel Caches

* microprocessor with multiple cores.



بخزن الـ I$ والـ D$ (private) → high miss rate ← لأنهم سريعين بس قليلين الدم
سوا لجيت هو combined فبكونوا مدعمين بـ Secondry cache (private)

→ large , slower

بفيدنا إنو الدم أكبر إنه بقلل الـ Miss rate فبزيد الـ Performance أكبر وأبطأ من L2 و Shared تخني كل الـ Processors
بستخدموه مش processor واحد ( Shared ) وهاد اشي مفيد لما الـ Processor
يتعاملوا مع نفس الـ data والـ code فبك بقلل الـ "Miss" وبزيد الـ "hit Rate" ... Rate"

SRAM

---

# Slide 73 : Dependability Measures :

(بالساعات أو الأيام بقاس ) من لما شغلناه لحتى تعطل ← ( بكون أكبر ما يمكن )



MTTF
ok
Fault

MTTR
MTBF → مجموع الزمنين

زمن التصليح ← ( كلما كان أصغر كل ما كان أفضل )

Slide 73 ؛ Improving Availabilty.

① Increase MTTF :

- Fault Avoidance : من البداية أتجنب الخطأ بإنه ال software و
ال hardware يكونوا أفضل ما يمكن ومن الصعب أحصل عليها لأنها غالية جدًا بتكون

- Fault tolerance : يعني لو تعطل عندي اشي يكون في إله بديل عشان ما توقف الشغل مثل إنه يكون فيه موتورين بتبريد مثلًا .

- Fault Forecasting : أتنبأ بالأعطال اللي ممكن تصير وبعمل على تبديل القطعة لما أعرف إنه حيخلص عمرهم أو حيخرب .

② Reduce MTTR : أحسن ال tools عشان تنبهوني بالأخطاء وتستخدموا وبساعدني أكتشف العطل و الآن لنقلل زمن التصليح : 1- يكون عندي مكان أدخل عليه العطل أحله بين ما أأثر عالباقي 2- نزيد أعداد المصلحين 3- نستعمل harddisck ثاني فلازم يكون عندي قطع غيار من أول احتياطًا .

Slide 74 ؛ Hamming Distance :

Distance = 1

Example Binary

| 0 | 00 |
| 1 | 01 |
| 2 | 10 |
| 3 | 11 |



Distance = 2

Parity

Ex: 011 → 0 ✓
010 → 1 ✗

Example Binary with Parity

| 0 | 00 | 0 |
| 1 | 01 | 1 |
| 2 | 10 | 1 |
| 3 | 11 | 0 |

even



وهون نفس ال
المبدأ من 0 ل 1 بنحتاج نعمل تغييرين هي أول تغيير وثاني تغيير نغير ال Parity
لعدد ال 1 زوجي بنحط 0 ولو وردي بنحط 1 من 0 ⟵ 1

## Distance = 3



‏بـ distance = 2 ← لو اختار وحدة من الحمر حيتحير لأي سودة يروح بس

‏بـ distance = 3 ← لو اختار وحدة من الحمر بخنار السودة الأقرب إلك

## Slide 76 : Encoding SEC :

Example
Data

1001 1010

| Bit position | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Encoded data bits | | P1 | P2 | d1 | P4 | d2 | d3 | d4 | P8 | d5 | d6 | d7 | d8 |
| Parity bit Coverate | P1 | X | | X | | X | | X | | X | | X | |
| | P2 | | X | X | | | X | X | | | X | X | |
| | P4 | | | | X | X | X | X | | | | | X |
| | P8 | | | | | | | | X | X | X | X | X |

0 1 1 1 0 0 1 0 1 0 1 0

From cpu to Memory   عددنا 1 زوجي دا   ‏هنا اللي بتخزن هنا

From Memory to cpu   0111 0010 1110          P8 P4 P2 P1

Single Error ‏في اختلاف عن اللي حليناه فعني            1  0  1  0

‏بين الأزرق والأحمر

‏لما حسبت الـ parity بين، إنه في error لأنه هو 0 والرقم الي طلعلي الـ = 10

‏وهو المكان اللي حصل فيه الخطأ والحل، إني أقلب قيمته للقيمة الصحيحة

⟹ 0111 0010 1110

0 111 0 110 1010

## Slide 77 :
Distance = 4



‏من Valid و Valid ‏بنا 4 changes  ما دنعرف مين زوج المسافة
‏code    code    منشاوية فبسهج نعمل
SEC / DED

## Slide 79 & RAID

### RAID 0

Data ↘  ↗ Disck

| A1 | A2 |
|----|----|
| A3 | A4 |
| A5 | A6 |
| A7 | A8 |

لوهاد تعطل فبنتكين
ضاعت كل ال Data
اللي عليه

### RAID 1

| A1 | A1 |
|----|----|
| A2 | A2 |
| A3 | A3 |
| A4 | A4 |

* لازم يكون عدد زوجي من ال disck
* لو تعطل واحد بقرأ من ال Mirror ثانية
* واللي تعطل بستبدله و بنقل عليه ال data

### RAID 2

| A1 | A2 | A3 | A4 | $AP_1$ | $AP_2$ | $AP_3$ |
|----|----|----|----|----|----|----|
| B1 | B2 | B3 | B4 | $BP_1$ | $BP_2$ | $BP_3$ |
| C1 | C2 | C3 | C4 | $CP_1$ | $CP_2$ | $CP_3$ |
| D1 | D2 | D3 | D4 | $DP_1$ | $DP_2$ | $DP_3$ |

Disck0  Disck1  Disck2  Disck3  Disck4  Disck5  Disck6

* غالي وبسبب ال hardware

### RAID 3

| A1 | A2 | A3 | Ap (1-3) |
|----|----|----|----------|
| A4 | A5 | A6 | Ap (4-6) |
| B1 | B2 | B3 | Bp (1-3) |
| B4 | B5 | B6 | Bp (4-6) |

N = 3

## RAID 4

| A1 | A2 | A3 | Ap |
|----|----|----|----|
| B1 | B2 | B3 | Bp |
| C1 | C2 | C3 | Cp |
| D1 | D2 | D3 | Dp |

| A | B | XOR |
|---|---|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

إذا في تغيير لح نعكس الـ Parity

و بنستخدم ناتج الـ Xor مع الـ Parity القديمة ونعملهم Xor بطلعوا الـ Parity الجديدة .

# RAID 5: Distributed Parity

( أفضل واحد في الوقت الحاضر )

## N + 1 disks

- Like RAID 4, but parity blocks distributed across disks
  - Avoids parity disk being a bottleneck

## Widely used

صاروا كلهم زي بعض والباريتي موزع



RAID 4                    RAID 5

---

⑦

# RAID 6: P + Q Redundancy

كشان لو بالصفة تعطل بكعنك 2 ماتو وحلينا وليقنا قادر يرجع الداتا

## N + 2 disks

- Like RAID 5, but two lots of parity
- Greater fault tolerance through more redundancy

## Multiple RAID → للمؤسسات الكبيرة اللي الها أكثر

- More advanced systems give similar fault tolerance with better performance
- Example RAID 51

# RAID Summary

RAID can improve performance and
availability كانه بدال ما يكون عندك Disck واحد يكون اكثر
من Disck

- High availability requires hot swapping → إذا تعطل اشي ما يوقف الجهاز دببدل المعطل بجديد

Assumes independent disk failures → إذا تعطل واحد البقية ما بتعطلوا

- Too bad if the building burns down! إذا احترقت البناية و داح كل ال RAID والاحل تروح لل Multiple RAID مثل RAID 51

# Contents

# Virtual Machines

يستخدم مثلاً لو بدي أشغل برنامج ال operating system، المختلف، اللي عندي

Host computer emulates guest operating system
and machine resources

ليقلد أو يحاكي

ما يشغل كل ال computer، بس ال cores، اللي حطت عليهم

- Improved isolation of multiple guests
- Avoids security and reliability problems
- Aids sharing of resources → ياخذ من ال resources، ينظم كل App قديه

كل VM، ما تشوف إلا ملفاتها ولو تعطلت، ما بتخرب غيرها

Virtualization has some performance impact

ممكن
- Feasible with modern high-performance comptuers

Examples

① ▪ IBM VM/370 (1970s technology!)
② ▪ VMWare   oracle   هي شركة
③ ▪ Microsoft Virtual PC

---

# Virtual Machines

cloud
↓

لو خرب واحد من ال App، ما يتأثروا الباقي

على الجهاز
↓



App A, App B, App C, App D, App E, App F

Host Operating System (EX: windows10)

Infrastructure (HW)

Virtual Machine — App A — Guest Operating System

Virtual Machine — App B — Guest Operating System

Virtual Machine — App C — Guest Operating System

Hypervisor

Infrastructure

نوع من أنواع ال ده بخطك تحمل أكثر من VM

# Virtual Machine Monitor

Maps virtual resources *(Guest)* to physical resources

- Memory, I/O devices, CPUs

Guest code runs on native machine in user mode → *User Mode مشتغل بال*

- 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

# Instruction Set Support

- User and System modes
- Privileged instructions only available in system mode *اللي بيتعاملوا مع الهاردوير وبيعملوا د input/output*
  - Trap to system if executed in user mode
- All physical resources only accessible using privileged instructions
  - Including page tables, interrupt controls, I/O registers

# Contents

. . .

*باستخراجها بحبر* ←
*ال access بطيئ*
*فبنلجأ لـ*

*ما في Process بنقدر ندخل* ← Memory Protection

*بـ code أو data لـ Process تانب وهذا فائدة الـ Virtual Memory*

*موجودة بكل أنظمة الحواسيب ،*
*بيني البرامج بنفرض اشي والحقيقة* ← *ذاكرة افتراضية*
*شي اخر*

# Virtual Memory

- Use main memory as a "cache" for secondary (disk) storage

*يعني الضا* ← *hardware*
*الفاي بنخل الـ Data*
*من الـ Memory*

- Managed jointly by CPU hardware and the operating system (OS)

MM

*الـ Cache*
*لكن نقل الـ Data*
*بين الـ Memory*
*والـ disk بنم*
*بتعاون*
*Software+*
*hardware*

- Programs share main memory → P₁ P₂ P₃
  - Each gets a private virtual address space holding its frequently used code and data
  - Protected from other programs
- CPU and OS translate virtual addresses to physical addresses
  - VM "block" is called a page
  - VM translation "miss" is called a page fault

Scanned with CamScanner

# Sharing the Physical Memory



Figure 3

Process 1
Virtual Address
Space (4 GB)

Process 2
Virtual Address
Space (4 GB)

Process 3
Virtual Address
Space (4 GB)

Physical RAM
(12GB example)

PAE Ext.

4GB
x86
Limit

# Address Translation

- Fixed-size pages (e.g., 4K)



Process
Virtual addresses

DRAM
Main Memory
Physical addresses

Address translation

Disk addresses

من cpu
Virtual address

47 46 45 44 43 ............... 15 14 13 12 11 10 9 8 ....... 3 2 1 0

| Virtual page number | Page offset |

Translation

39 38 37 ............... 15 14 13 12 11 10 9 8 ....... 3 2 1 0

| Physical page number | Page offset |

Physical address

بيتمكنف ← اختار الـ Byte اللي بدي ياها من الـ 4KB اللي موجودين بالـ Page ولأنهم نفس الترتيب فالـ Page offset مافي داعي نعملها translation

رقم الـ page ← فلو إجابي الـ Page لازم أعملها translation لإيفا موجودة هنالك الـ Page 2 الـ OS حطواهنالك

\* الـ Virtual Addresses ممكن يكونوا أكبر من الـ Physical Addresses فمش كلهم إلهم مكان بالـ MM فبكون في أشياء بالـ Disk

هو اللي بيروح الـ Mem الـ access ويستخدمه ليروح على الـ Mem ولو ماكانت الـ Page الـ Mem بيحصل Page fault

# Page Fault Penalty

← أبطأ هنة أ لغامرة من ال M فلو الشخلة ما قنناها بال Memory وصبر ناح

Page fault واضطرينا

نرح لل Disk فيكون

Penalty نس

## On page fault, the page must be fetched from disk →

بديل ال Disk هي الأيام ؟
SSD → Flash ال معمولمن

لأنه يقاس بال M sec

- Takes millions of clock cycles → M sec

عليتصمقدرة ال hardware ما يقدر

بيطلعها ال os هو الي بيطلعها لما يحبر في Page fault

- Handled by OS code →

Data اللي مش موجودة بال MM

## Try to minimize page fault rate →

ال os مكتوبة بإزوا نقلل قد ما نقدر Page fault ال

بندخل وبقوف المتنكلة و بجيب ال

① - Fully associative placement

② - Smart replacement algorithms

يعني إذا جيت اشوي من ال disk بحطوبأي مكان بدي يلي يكون فاضي

Least Recently used / Random مثلة

أحسن بس أقدر

معتمد على ال hardware بس لل Software بنقدر نكتب

---

# Page Tables →

ال os حتى يقدر يعمل translation بيكون عنده Page tables

## Stores placement information

→ Virtual Page وين موجودة بال Main M

يعني بخزن كل

PT
3
2
1 → PTE
0

- Array of page table entries, indexed by virtual page number →

عنحريق ال VPN بنوصل لل Array و بنطلع ال PPN ( Physical Page Number )

لكل Process

- Page table register in CPU points to page table in physical memory

الـ os يعملوا PT ولما process 1 تكون شغالة بنستخدم ال PT تاعوا

PTR فا ال باشر كبراية الـ PT الـ Process اللي شغال

## If page is present in memory (Page hit)

- PTE stores the physical page number
  PPN

- Plus other status bits (referenced, dirty, …)

## If page is not present

- PTE can refer to location in swap space on disk

# Translation Using a Page Table

Page table register ← CPU

48 bit = Virtual address ← CPU بيجي من ال CPU

47 46 45 44 43 ················ 15 14 13 12 11 10 9 8 ······ 3 2 1 0

| Virtual page number | Page offset |
|---|---|

36

Valid    Physical page number

12 = log 4 KB = 12

Page table

لم يستعمله لطول ال Virtual Address

Physical address

If 0 then page is not present in memory

28

39 38 37 ···················· 15 14 13 12 11 10 9 8 ······ 3 2 1 0

| Physical page number | Page offset |
|---|---|

Physical address →

اللي بستخدمه Mem ال Access

باشر كتابة ال RAM

لو كانت ا بكون Page fault وبكون موجود بال disk

**Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 98**

# Mapping Pages to Storage

Virtual page number

Page table
Physical page or
Valid    disk address

| |
|---|
| 1 |
| 1 |
| 1 |
| 1 |
| 0 |
| 1 |
| 1 |
| 0 |
| 1 |
| 1 |
| 0 |
| 1 |

Physical memory

Disk storage

# Replacement and Writes

- To reduce page fault rate, prefer least-recently used (LRU) replacement
  - Reference bit (aka use bit) in PTE set to 1 on access to page
  - Periodically cleared to 0 by OS ← كل لحظ ثانية الـ os بصحح الـ use bit
  - A page with reference bit = 0 has not been used recently
- Disk writes take millions of cycles
  - Block at once, not individual locations ← ننكتب بالـ Mem ما بنبعتوا للـ disk
  - Write through is impractical ← مش عملي
  - Use write-back ← هو العملي
  - Dirty bit in PTE set when page is written

MK

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 100

# Fast Translation Using a TLB
← طرق نسرع الـ translation

- Address translation would appear to require extra memory references $VA \rightarrow PT \rightarrow PA \rightarrow M$
  - One to access the PTE
  - Then the actual memory access ← أبطأ جزءًا من الـ Process
- But access to page tables has good locality ← الـ Page الوحدة بنستخدموا مرات عديدة
  - So use a fast cache of PTEs within the CPU
  - Called a Translation Look-aside Buffer (TLB) ← دي نكتب الـ TLB
  - Typical: 16–512 PTEs, 0.5–1 cycle for hit, 10–100 cycles for miss, 0.01%–1% miss rate ← good locality
  - Misses could be handled by hardware or software

لـ الـ cpu بنروح لـ TLB وهي ماخية كل شي في جزء من الـ translations ولو ماالقيت جبنا بنك نروح لـ PT وحتى نروح PT وتجيب المطلوب ونحطه بالـ TLB hardware بنستخدم بعض الـ System او software

MK

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 1

Scanned with CamScanner

# Fast Translation Using a TLB

**CPU**

Virtual page number

TLB →

Page Table الـ cache

يحمل بعض الـ translationsالـ

**Physical memory**

كلمتكو

Fully or set assosiative

**Page table**

**Disk storage**

# TLB Misses

① **If page is in memory** → الـ Page موجودة بالـ Memory

لسا مش موجودة بالـ TLB الـ PT

- **Load the PTE from memory and retry**
- **Could be handled in hardware** → يعتمد فيه الـ PT

Simple فنستخدم الـ hardware

- Can get complex for more complicated page table structures
- **Or in software**
  - Raise a special exception, with optimized handler

يعني الـ Instructions محسوبة بدقة

② **If page is not in memory (page fault)** → System الـ

- **OS handles fetching the page and updating the page table**

بكل الـ بتصرف الـ OS os

- **Then restart the faulting instruction**

# TLB Miss Handler

TLB miss indicates

① ■ Page present, but PTE not in TLB  هاي أبسط

② ■ Page not preset ( Page fault )

Must recognize TLB miss before
destination register overwritten → ld x, مثلا لو عندي

■ Raise exception  لو عملت TLB Miss مش لازم نخليها تكمل ولا تكتب
ع x, لإن الـ translation ما صار لسا .

Handler copies PTE from memory to TLB

■ Then restarts instruction

■ If page not present, page fault will occur

Software لأ الـ يتصرف

# Page Fault Handler → جزء من الـ os
بشتغل لما يحصر
page fault لنا

Use faulting virtual address to find PTE

Locate page on disk

Choose page to replace

■ If dirty, write to disk first

Read page into memory and update page
table

Make process runnable again

■ Restart from faulting instruction

Ex:  ld  — TLB — hit ✓
              — Miss — PT — Valid ✓
                            — Page fault

Scanned with CamScanner

# TLB and Cache Interaction



**Virtual address**

```
       ...... 14 13 12 11 10 9 ........ 3 2 1 0
31 30 29
Virtual page number                 Page offset
```

Valid Dirty · Tag · Physical page number

Physical page number · Page offset

—Physical address—

Physical address tag · Cache index · Block offset · Byte offset

PA

Valid · Tag · Data

Data

If cache tag uses physical address

- Need to translate before cache lookup

Alternative: use virtual address tag أسرع

- Complications due to aliasing
  - Different virtual addresses for shared physical address

لا يعني ممكن يكون عندك
2P ل shared وهي one page
بطلعولها different Addresses

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 106

## RAID5

| A₁ | | A₂ | | A₃ | | Aₚ |
|----|----|----|----|----|----|----|
| B₁ | | B₂ | | Bₚ | | B₃ |
| C₁ | | Cₚ | | C₂ | | C₃ |
| Dₚ | | D₁ | | D₂ | | D₃ |

## RAID 6

| A₁ | | A₂ | | A₃ | | Aₚ | | A𝓺 |
|----|----|----|----|----|----|----|----|----|
| B₁ | | B₂ | | Bₚ | | B𝓺 | | B₃ |
| C₁ | | Cₚ | | C𝓺 | | C₂ | | C₃ |
| Dₚ | | D𝓺 | | D₁ | | D₂ | | D₃ |
| E𝓺 | | E₁ | | E₂ | | E₃ | | Eₚ |

* الـ P يحسب سطريًا أفقيًا.

* الـ Q تحسب بطريقة أخرى كسطر مائل مش أفقي (قطري) diagonal

| RAID 5 | | | RAID 51 RAID1 | | | RAID 5 | | |
|--------|--|--|---------------|--|--|--------|--|--|
| A1 | | B1 | | P1 | | A1 | | B1 | | P1 |
| A2 | | P2 | | B2 | | A2 | | P2 | | B2 |
| P3 | | A3 | | B3 | | P3 | | A3 | | B3 |
| B4 | | A4 | | P4 | | B4 | | A4 | | P4 |

Slide 94 : Sharing the physical Memory :



Physical Memory

* مقسمة إلى Pages و الـ page الوحدة أكثر شي = 4kB .
* بالأخص كل Process الـ Pages خاصة فيها, الا إذا بدنا نعمل Sharing زي رقم 1 .
* الأبيض يعبر عن Pages فاضية .
* المبرمج بفترض إنه بيقدر يستخدم أي Address من أكبر رقم لأقل رقم .
* مسؤولية الـ cpu hardware + os بيوزعها الـ Virtual spaces بالـ Physical Memory ( Pages )

# Slide 106 : TLB- Cache Interaction :

## * Physical Address Tag :

* يعني الـ cache يكون فيه physical Address مهو ـ عقة الـ Memory
فلازم أول يكون الـ TLB وهيك بكلانكيث ومحو الـ Cache أكبر .

| P | VA → | TLB | PA → | $ | → | M | |
|---|---|---|---|---|---|---|---|

## * Virtual Address Tag :

* بنخلي الـ tags الي في الـ Cache كأنها تكون Virtual Address ومعناه ما في حاجة
نعمل Translation الا وقت الـ Miss .

| P | VA → | $ | VA → | TLB | PA → | M | |
|---|---|---|---|---|---|---|---|

## * Physical Address Tag : cache Size = Page Size * associativity
لحيب يكون الحالة الخاصة .

# Memory Protection

Parallel Program يشوفوا بعض

Different tasks can share parts of their (read) (write) (Fetch)
r    w     F

virtual address spaces

access في حاجة معينة
- But need to protect against errant access

هو اللي يكتشف الأخطاء
- Requires OS assistance

هو اللي ينسق ويحدد
Hardware support for OS protection

r w F
$110 \rightarrow$ R W
$100 \rightarrow$ R W
$001 \rightarrow$ read only / code
$000 \rightarrow$ OS

- Privileged supervisor mode (aka kernel mode)

- Privileged instructions

- Page tables and other state information only
  accessible in supervisor mode → بنودنيا الـ operating system

- System call exception (e.g., ecall in RISC-V)

# Contents

# The Memory Hierarchy

## The BIG Picture

- Common principles apply at all levels of the memory hierarchy
  * كل اشي هو cache الي تحته موجود فيه جزء من معلومات اللي تحته .
  - Based on notions of caching
- At each level in the hierarchy
  - Block placement → إذا جبنا اشي من تحت لفوق وين نحطه .
  - Finding a block → وين هذا الـ address إذا موجود بال level اللي احنا فيه .
  - Replacement on a miss ← إذا كل الأماكن مليانة وين نحط .
  - Write policy → لما بدنا نكتب كيف نكتب

# Block Placement

- Determined by associativity
  - Direct mapped (1-way associative)
    - One choice for placement      كل Mem block اله مكان واحد بتحدد جاد index .
  - n-way set associative (2/4/8/16)هيك مش اكثرمن
    - n choices within a set
  - Fully associative
    - Any location      ماانا index كل الأماكن بال Cache متاحين .
- Higher associativity reduces miss rate
  - Increases complexity, cost, and access time

مشكلة higher associativity

# Finding a Block

| Associativity | Location method | Tag comparisons |
|---|---|---|
| Direct mapped | Index → block ال بعتمد عددد بال Cache | 1  (=) |
| n-way set associative | Set index, then search entries within the set | n |
| Fully associative | Search all entries  $ | #entries |
| | Full lookup table  VM | 0 |

في فيها Page table

- Hardware caches
  - Reduce comparisons to reduce cost    VA → [PT] → PA
- Virtual memory
  - Full table lookup makes full associativity feasible
  - Benefit in reduced miss rate

Scanned with CamScanner

# Replacement

- Choice of entry to replace on a miss
    - Least recently used (LRU)
        - Complex and costly hardware for high associativity
    - Random
        - Close to LRU, easier to implement
- Virtual memory
    - LRU approximation with hardware support

recently used : يقسم ال Page لجزئين
Replacement ال مع و بينم ال not recently used و
not recently used من ال
وال bit اللي بال VM اللي بستخدمها ال os
هي ( ١ bit ) ← يعني ال os كل
فترة وفترة بحط Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 112

بزودها o بلف عكل ال Pages وبحط حوا o بعين
ال Page اللي بصير لها access بتحول د ١ وهيك بنكون

أشرنا عال | Pages اللي |

# Write Policy

- Write-through
    - Update both upper and lower levels
    - Simplifies replacement, but may require write buffer
- Write-back
    - Update upper level only
    - Update lower level when block is replaced
    - Need to keep more state
- Virtual memory
    - Only write-back is feasible, given disk write latency

Scanned with CamScanner

# Sources of Misses

- **Compulsory misses (aka cold start misses)** إجباري

  لما أول مرة بدنا نفعل و نشغل البرنامج لأ
  - First access to a block
  الـ Caches بكونوا لساهاضيين فأول access اجباري تكون Miss لأنها مشاهوجدة.

- **Capacity misses**
  لحجم ما بكون واسع في الـ cache فلو عنلنا →
  - Due to finite cache size معلومة ما لقيتلها لأنه ادخط مكانها الشيء الثاني access
  - A replaced block is later accessed again فبصير في Miss
  إذا بدي أقلل الـ capacity Miss بزيد حجم الـ cache

- **Conflict misses (aka collision misses)**
  - In a non-fully associative cache
  - Due to competition for entries in a set
  - Would not occur in a fully associative cache of the same total size

يعني يصرف بالـ cache مكان معين، الـ index معين، انه الـ P بطلب من الـ M عددها blocks



أكثر من واحد بالوم

نفس الـ x اعطها بالـ cache فلو طلب الأول بيجي عادي لو طلب الثاني بده يجيبه الـ cache بس
يطرد الأول فبصير الأول والثاني يضلهم يطردوا بوض.

# Cache Design Trade-offs

لأ مع انه الـ cache في أماكن فاضية بس اشروا عنفس الـ index

B1 →
B2 →

$$* AMAT = HT + \text{Miss rate} * \text{Miss Penalty}$$

| Design change | Effect on miss rate | Negative performance effect |
|---|---|---|
| Increase cache size | Decrease capacity misses | May increase access time |
| Increase associativity | Decrease conflict misses | May increase access time |
| Increase block size | Decrease compulsory misses | Increases miss penalty. For very large block size, may increase miss rate due to pollution. |

لما الـ block size بكبر بصير More data كثر Miss أجيب أقل فلو بحتاجها بكون وقت عحالي Misses بالمستقبل

# Data Cache Miss Rate



**Miss rate** (نسبة كاذبة الأخطاء)

15%
12%   1 KiB
9%    2 KiB
6%    4 KiB
      8 KiB
3%    16 KiB
      32 KiB      64 KiB      128 KiB
0

One-way    Two-way    Four-way    Eight-way

**Associativity**

\* العلاقة بين الـ Miss rate والـ Associativity و بتقل الـ Miss rate كل ما يزيد الـ Associativity
السبب إنه الـ conflict Missis بتقل

\* لما يكبر الـ cache بتقل الـ Miss rate وبعدها الـ capacity Missis

# Contents

بتُقل \* لما نوصل لـ 128 KiB ولو زدنا الـ Associativity بتفل ثابت الـ Miss rate والسبب وجود الـ compulsory Missis فلا يمكن ينقص أكثر من هيك

# Cache Control

Example cache characteristics

dirty bit ~~فيه~~

- Direct-mapped, write-back, write allocate
- Block size: 4 words (16 bytes)
- Cache size: 16 KB (1024 blocks)
- 32-bit byte addresses
- Valid bit and dirty bit per block
- Blocking cache → write_back ~~لأنه~~
  - CPU waits until access is complete

| 31 | 14 13 | 4 3 | 0 |
|---|---|---|---|
| Tag | Index | Offset | |

18 bits — 10 bits — 4 bits

↳ log 1K = 10    ↳ lg 16 = 4

32 - 10 - 4 = 18

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 118

# Interface Signals

Read/write (Mem)
Valid (Mem)
Ready (CPU)

cache bus/back side bus    FSM    Memory bus / Front side bus

cache و cpu بين    عن طريقه بتجيب الداتا من الـ M الى الـ Cache

الى الـ Cache

Miss لـ



CPU — Read/Write, Valid, Address 32, Write Data 32, Read Data 32, Ready — Cache — Read/Write, Valid 0→1, Address 32, Write Data 128, Read Data 128, Ready — Memory

بكتب ه 1 نوأ

بنحير 1 عشان

لا يجيب لا Cache

على الغلط

كبير للتسهيل

لو بدو ينقول الـ

لو بدو ينقول الـ cpu هو هل Ready و لا لا
1 = Ready و لكن لو كان Hit يكون
Ready و لو كان Miss بياخذ زمن ليجهز Ready

Multiple cycles per access

Miss

Chapter 5 — Large and Fast: Exploiting Memory Hierarchy — 119

# Finite State Machines (sequential)

- Use an FSM to sequence control steps
- Set of states, transition on each clock edge
  - State values are binary encoded
  - Current state stored in a register
  - Next state
    $= f_n$ (current state, current inputs)
- Control output signals
  $= f_o$ (current state)

| Sn | I/P | State | o/P |
|----|-----|-------|-----|
| :  | :   | :     | :   |

*خطوات تحكم متسلسلة* & Sequential circuits



Combinational control logic
Datapath control outputs
Outputs
Inputs
Inputs from cache datapath
State register
Next state

Sequential circuits → input و الـ out بعتمد علي قيمة ذاكرة و الـ State
Combinational circuits → ما فيها ذاكرة و الـ out بتعتمد علي الـ input بس

# Cache Controller FSM

State البدائية
الـ cache بتكون بالـ Idle state



**S0 Idle**

Cache Hit
Mark Cache Ready

Valid CPU request

**S1 Compare Tag**
If Valid && Hit,
Set Valid, SetTag,
if Write Set Dirty

Could partition into separate states to reduce clock cycle time

Cache Miss and Old Block is Clean
D=0
ممكن مكانها نستبدل الـ block

Cache Miss and Old Block is Dirty
D=1 → block يعني الـ قبل ما أعمل Allocate لازم أعمل منه write-Back

أعقد حالة

Memory Ready

**S2 Allocate**
Read new block from Memory

Memory not Ready

Memory Ready

**S3 Write-Back**
Write Old Block to Memory

Memory not Ready

# Contents

MK

هتر ابطة ومتعاضرة ومتزاحمة

# Cache Coherence Problem

- Suppose two CPU cores share a physical address space
  - Write-through caches → Write-through بالـ هاي المشكلة بنحير بالـ
    والـ wnte-back

| Time step | Event | CPU A's cache | CPU B's cache | Memory |
|-----------|-------|---------------|---------------|--------|
| 0 | | | | $X = 0$ |
| 1 | CPU A reads X | 0 | | 0 |
| 2 | CPU B reads X | 0 | 0 | 0 |
| 3 | CPU A writes 1 to X | 1 | 0 | 1 |

رح يكتب عالـ 1 = Memory والـ cache

وهون في مشكلة لانه الـ Caches صار فياينهم نسختين من X فلو B cpu قرأ X حيقرلها مش 1 وهاي مشكلة والحل في الحلف

# Coherence Defined

- Informally: Reads return most recently written value

  تعريف الـ cache coherence — قيمة

  يعني لما الـ cpu تعمل read لازم تقرأ أجدد

- Formally:

  التعريف المطول

  ① P writes X; P reads X (no intervening writes) ⇒ read returns written value

  بعد فترة زمنية معقولة

  ② $P_1$ writes X; $P_2$ reads X (sufficiently later) ⇒ read returns written value
     - c.f. CPU B reading X after step 3 in example

  ③ $P_1$ writes X, $P_2$ writes X ⇒ all processors see writes in the same order
     - End up with the same final value for X

  نشوف بأي ترتيب كتبوا و بناخذ آخر كتابة كتبوها أو انكتبت

# Cache Coherence Protocols

- Operations performed by caches in multiprocessors to ensure coherence
  - Migration of data to local caches
    - Reduces bandwidth for shared memory
  - Replication of read-shared data
    - Reduces contention for access

  أنواع الـ cache coherence ؟

- ① Snooping protocols
  - Each cache monitors bus reads/writes
- ② Directory-based protocols
  - Caches and memory record sharing status of blocks in a directory

  دليل

# Invalidating Snooping Protocols

- Cache gets exclusive access to a block when it is to be written
  - Broadcasts an invalidate message on the bus
  - Subsequent read in another cache misses
    - Owning cache supplies updated value

| CPU activity | Bus activity | CPU A's cache | CPU B's cache | Memory |
|---|---|---|---|---|
| | | I | I | X = 0 |
| CPU A reads X | Cache miss for X | 0 | I | 0 |
| CPU B reads X | Cache miss for X | S 0 | S 0 | 0 |
| CPU A writes 1 to X | Invalidate for X | M 1 | I | 0 |
| CPU B read X | Cache miss for X | S 1 | S 1 | 1 |

# Memory Consistency

- When are writes seen by other processors
  - "Seen" means a read returns the written value
  - Can't be instantaneously → بسبب إنه في تبليد بين Processorال
- Assumptions
  - A write completes only when all processors have seen it
  - A processor does not reorder writes with other accesses كما في sequential consistency
    هاي الأيام إنه بخلي البرامج بطيئة وإذا موجود لازم يكون ماسكه
- Consequence
  - P writes X then writes Y
    ⇒ all processors that see new Y also see new X
  - Processors can reorder reads, but not writes

# Contents

# Multilevel On-Chip Caches

*[handwritten: 1 GHz]*   *[3-4 GHz]*   *[P, I$ D$, 32 k  32 K]*

*[الحفظة]*

| Characteristic | ARM Cortex-A53 | Intel Core i7 |
|---|---|---|
| L1 cache organization | Split instruction and data caches | Split instruction and data caches |
| L1 cache size | Configurable 16 to 64 KiB each for instructions/data | 32 KiB each for instructions/data per core |
| L1 cache associativity | Two-way (I), four-way (D) set associative *[ب]* *[المقطع كله أكبر]* | Four-way (I), eight-way (D) set associative |
| L1 replacement | Random | Approximated LRU |
| L1 block size | 64 bytes *[Hit  Miss]* | 64 bytes *[64/4=16 word = 8 Dw each block]* |
| L1 write policy | Write-back, variable allocation policies (default is Write-allocate) | Write-back, No-write-allocate *[Hit   Miss]* |
| L1 hit time (load-use) | Two clock cycles | Four clock cycles, pipelined *[اسرع]* |
| L2 cache organization | Unified (instruction and data) | Unified (instruction and data) per core |
| L2 cache size | 128 KiB to 2 MiB | 256 KiB (0.25 MiB) |
| L2 cache associativity | 16-way set associative | 8-way set associative |
| L2 replacement | Approximated LRU | Approximated LRU |
| L2 block size | 64 bytes | 64 bytes |
| L2 write policy | Write-back, Write-allocate | Write-back, Write-allocate |
| L2 hit time | 12 clock cycles *[أكبر ودقة]* | 10 clock cycles *[الحجم واسع]* |
| L3 cache organization | – | Unified (instruction and data) |
| L3 cache size | – | 8 MiB, shared |
| L3 cache associativity | – | 16-way set associative |
| L3 replacement | – | Approximated LRU |
| L3 block size | – | 64 bytes |
| L3 write policy | – | Write-back, Write-allocate |
| L3 hit time | – | 35 clock cycles |

# 2-Level TLB Organization

| Characteristic | ARM Cortex-A53 | Intel Core I7 |
|---|---|---|
| Virtual address | 48 bits | 48 bits |
| Physical address | 40 bits | 44 bits |
| Page size | Variable: 4, 16, 64 KiB, 1, 2 MiB, 1 GiB | Variable: 4 KiB, 2/4 MiB |
| TLB organization | 1 TLB for instructions and 1 TLB for data per core<br><br>Both micro TLBs are fully associative, with 10 entries, round robin replacement<br>64-entry, four-way set-associative TLBs<br><br>TLB misses handled in hardware | 1 TLB for instructions and 1 TLB for data per core<br><br>Both L1 TLBs are four-way set associative, LRU replacement<br><br>L1 I-TLB has 128 entries for small pages, seven per thread for large pages<br><br>L1 D-TLB has 64 entries for small pages, 32 for large pages<br><br>The L2 TLB is four-way set associative, LRU replacement<br><br>The L2 TLB has 512 entries<br><br>TLB misses handled in hardware |

*Handwritten annotations:*
- TLB organization → Cache, PTE, VA يحول, PA إلى
- The micro TLBs are L1 TLBs that are backed by two 64-entry L2 TLBs.
- $2^{44}$ = 16 TB
- 128+7   64+32
- L1   TLB I   TLB D
- L2   TLB : 512

# Supporting Multiple Issue

- Both have multi-banked caches that allow multiple accesses per cycle assuming no bank conflicts

- Other optimizations
  - Return requested word first
  - Non-blocking cache
    - Hit under miss
    - Miss under miss

# Slide 123 : Cache Coherence.

PA    PB

$ 0→1    0 $  →

يظل فيه القيمة القديمة لـ X القيمي

Bus

M X=0 →1

فيصير بطئ في تعاون بين A و B
وهذا هو cache coherence problem
أي شي يروح عالـ M بين الـ Bus

| Time step | Event | CPU A's cache | CPU B's Cache | Memory |
|---|---|---|---|---|
| 0 | | | | 0 |
| 1 | cpu A reads X | 0 | | 0 |
| 2 | cpu B reads X | 0 | 0 | 0 |
| 3 | cpu A writes I to X | I | 0 | I |

# Slide 126 : Snooping Protocols.

PA    PB

$ 0→1    0 $
write Back

inv. ↓

Signal → تعني إنه
Exec. copy    حياخذ

M X=0    read Miss ← Bus

لما يراقب الـ cache الـ Bus ويشوف
صار فيه Exec. copy بخليه
invalid يعني بيشيل الرقم
والـ M خلا فيوا القيمة القديمة ولما رجعنا عملنا
read لـ B بكون Miss فيطلع عالـ Bus read Miss
وبياخذ الـ 1 من الـ A

| cpu activity | Bus Activity | CPU A's cache | CPU B's cache | Memory |
|---|---|---|---|---|
| | | | | 0 |
| cpu A reads X | Cache Miss for X | 0 | | 0 |
| cpu B reads X | Cache Miss for X | 0 | 0 | 0 |
| cpu A writes I to X | Invalidate for X | I | | 0 |
| cpu B reads X | Cache Miss for X | I | I | I |

Scanned with CamScanner

## Slide 127 : Memory Consistency.

| Processor A | Processor B | |
|---|---|---|
| Flag = False ; | | |
| | while (flag == false); | |
| X = Compute(); | | |
| | Print (X) ; | |
| Flag = True ; | | |

2 Processors بنعاملوا مع بعض : A بيعمل حسبة معينة و بيكتب قيمة جديدة لـ X
وبعدين flag= true يعني إنه خلّص.

B بدو يطبع القيمة الجديدة لـ X بس قبل ما يقراها ويطبعها بكون براقب الـ flag لو كان
flase بضلها بالـ while لو صارت الـ Flag= true بطلع من الـ Loop وبقدر إنه ينفذ
الـ Print . ولو في Sequential Consistency فلما تطلع الـ flag=t معناها كتب
قيمة X الجديدة فيشتغل البرنامج صح .

هاي الأيام ما في Sequential Consistency والـ Processors بتكتب وفي oot of order
execution فممكن إنه قيمة الـ Flag تظهر لـ B قبل ما تظهر له القيمة الجديدة لـ X
وبالتالي البرنامج ما بشتغل صح فالحل إنه المبرمجين هالأيام لما بكتبوا Parallel Programs
بعتمدوا على تقنيات خاصة ( For synchronisation ) حتى تضمنلهم إنه الشغل
صح.

# Supporting Multiple Issue

(لبس اجا خندكي بشكل عام)

→ ARM and Intel

- **Both** have multi-banked caches that allow multiple accesses per cycle assuming no bank conflicts (المخزن بتزجها)

- Other optimizations

  - ① Return requested word first

  - ② Non-blocking cache

    - Hit under miss

    - Miss under miss →

  - ③ Data prefetching

بساعدنا نقلل من هذا ال Capacity Missis

# Contents

# Pitfalls

- Byte vs. word addressing
  - Example: 32-byte direct-mapped cache, 4-byte blocks
    - Byte 36 maps to block 1
    - Word 36 maps to block 4
- Ignoring memory system effects when writing or generating code
  - Example: iterating over rows vs. columns of arrays
  - Large strides result in poor locality

# Pitfalls

- In multiprocessor with shared L2 or L3 cache
  - Less associativity than cores results in conflict misses
  - More cores $\Rightarrow$ need to increase associativity
- Using AMAT to evaluate performance of out-of-order processors
  - Ignores effect of non-blocked accesses
  - Instead, evaluate performance by simulation

# Pitfalls

- Extending address range using segments
  - E.g., Intel 80286
  - But a segment is not always big enough
  - Makes address arithmetic complicated
- Implementing a VMM on an ISA not designed for virtualization
  - E.g., non-privileged instructions accessing hardware resources
  - Either extend ISA, or require guest OS not to use problematic instructions

# Concluding Remarks

- Fast memories are small, large memories are slow
  - We really want fast, large memories ☹
  - Caching gives this illusion ☺
- Principle of locality
  - Programs use a small part of their memory space frequently
- Memory hierarchy
  - L1 cache ↔ L2 cache ↔ … ↔ DRAM memory ↔ disk
- Memory system design is critical for multiprocessors

# Slide 131 : Supporting Multiple Issue.



زوجي

Address 0 → Bank0 → Data 0

لوكان فردي بروح لهون ← والعكس

فردي

Address 1 → Bank1 → Data 1

"هذه الرسمة تتوضح توزيع cache على 2 Banks"
فبنجعل الـ Set اللي الـ Address تبعه زوجي تلاقيه
بـ Bank الـ 0 والـ Set اللي الـ Address الـ فردي تلاقيه
بـ Bank الآخر. ولما يكون في 2 Bank معناها
فيه 2 Ports فمعناه الـ P بعمل 2 access بنفس الوقت
وهيك بتكون الـ Performance

# Slide 131: Return requested word first.

• Early restart

| 1 | 2 | 3 | 4 |

\* Requested word : word 3

1 → 2 → 3 → 4

\* ☐ processing performed background

• Critical word first

| 1 | 2 | 3 | 4 |

3 → 4 → 1 → 2

\* الـ Block هون بتكون من 4 word في التصاميم القديمة كنا لازم نجيب w4,w3,w2,w1 بعدها نجيبهم كلهم نحطهم بالـ Cache و الـ P بياخد الـ Data من الـ cache بعد ما يوصل الـ Block كامل.

\* بالـ Early restart لما توصل الـ w3 ما بنستنى لتتخزن بالـ Cache ونعطيها للـ P اللي بدو ياها ديركت من الـ cache بلا ما يمروها لـ P.

\* بالـ Critical word first بدل ما ينجابوا w1, w2, w3, w4 بنطلبوا من الـ Memory من Address 3 فبيجوا w2, w1, w4, w3 فبنقلل الـ Miss penalty أكثر.

\* في الأجهزة الحديثة موجودة الطريقتين.

# Slide 131: Non-blocking cache.

cache miss
↓

Blocking Cache:

| CPU time | ↔No cache Access↔ | CPU time |

| Miss Penalty |

cache miss  hit  stall on use
↓ ↓ ↓

Hit Under Miss:

| CPU time | | CPU time |

| | Miss Penalty |

Cache Miss  Miss  Stall on use
↓ ↓ ↓

Miss under Miss :

| CPU time | | CPU time |

| Miss Penalty |

| Miss Penalty |

## Slide 131 : Data prefetching

| Access Block A |
|---|

لما الـ Processor يطلب الـ A الـ prefetch circuit بتستنى الأحداث ← وبتعمل prefetching الـ Next sequential Block وهذا اشي

| Prefetch Block A+1 |
|---|

ممتاز وبنوفر على حالة الـ Miss .

لـ الطريقة المعقدة الصحيحة :



**legend**

→ prefetch Request

→ prefetched Data

# Chapter 6

## Parallel Processors from Client to Cloud

لـ مثل الهواتف
والنوت بوك

*Adapted by Prof. Gheith Abandah*

# Contents

# Introduction

- Goal: connecting multiple computers to get higher performance
  - Multiprocessors ( بصير هيك اسم الجهاز )
  - Scalability, availability, power efficiency
- Task-level (process-level) parallelism
  - High throughput for independent jobs
- Parallel processing program
  - Single program run on multiple processors  Ex. SPMO
- Multicore microprocessors
  - Chips with multiple processors (cores)

*(handwritten annotations:)*
multiple processors are more Power efficiency than Single Processors
→ الـ multiple jobs ينفذها على multiple cores
Ex. SPMO ↓ Data
بحرف طاقة أقل ③  ② ① قدرات الـ Multip اكبر من Single الـ

# Hardware and Software

- Hardware
  - ① Serial: e.g., Pentium 4 كان موجود قبل 20 سنة
  - ② Parallel: e.g., quad-core Xeon e5345
  
  ممكن تنفذه ← لو تنفذهون بياخذ عمر واحد ويترك الباقي  بنفذ 4 program بنفس الوقت
- Software
  - ① Sequential: e.g., matrix multiplication
  - ② Concurrent: e.g., operating system معنا قسموه ← لأجزاء كل جزء اله برمجة خاصة فيه ويتعاونوامع بعض
  
  لو تنفذ ← هون رح يتنفذ بزمن أقل وبتنفذهوا وحدة ورا الثانية
- Sequential/concurrent software can run on serial/parallel hardware
  - Challenge: making effective use of parallel hardware

# What We've Already Covered

- §2.11: Parallelism and Instructions
  - Synchronization
- §3.6: Parallelism and Computer Arithmetic
  - Subword Parallelism → Performence بعطينا أعلى
- §4.10: Parallelism and Advanced Instruction-Level Parallelism
- §5.10: Parallelism and Memory Hierarchies
  - Cache Coherence

# Contents

MK

Chapter 6 — Parallel Processors from Client to Cloud — 6

# Contents

## 6.2 The Difficulty of Creating Parallel Programs

Parallel Programming

Amdahl's Law

Scaling

Strong and Weak Scaling

# Parallel Programming

- Parallel software is the problem
- Need to get significant performance improvement → البرنامج هش بس بيشتغل مع لا يشتغل مع و يكون الـ Performance عالي و التحسين جدًا عالي يكون .
  - Otherwise, just use a faster uniprocessor, since it's easier! لـ لو كان صعب عليك تعمل البرمجة بتحسين عالي فبنستخدم الـ Fster uniprocessor.
- Difficulties
  - Partitioning → تجزء الشغل اللي بدك تعمله لأجزاء
  - Coordination → خلال التنفيذ لازم يستقوا الـ cores كيف يبلشوا ويخلصوا مع بعض و يشتغلوا مزبوط
  - Communications overhead لما يحيروا الـ Processors يتعاونوا مع بعض بحيروا يحتاجوا بها وقت من بعض فبحير في زمن كبير غيرعن لما الـ Processor بيشتغل لحاله بدون ما يشارك Processors آخرين .

ولما نكتب Parallel Program لازم ننتبه لهاي الصعوبات

# Amdahl's Law

بحكيناعن الـ speed up اللي ممكن نحمله .

- Sequential part can limit speedup   (1-f) Sequential part   هو أحسن speed up ممكن نحطه بعنمه S
- Example: 100 processors, 90× speedup?
  - $T_{new} = T_{parallelizable}/100 + T_{sequential}$
  - $$Speedup = \frac{1}{(1-F_{paralleliable}) + F_{paralleliable}/100_P} = 90$$
  - Solving: $F_{parallelizable} = 0.999$
- Need sequential part to be 0.1% of original time وهنا فعلا صعب فهنا واحدمن الصعوبات اللي بتواجهنا

# Scaling Example

أمثلة بنوجبنا العلاقة بين حجم
المسألة وعدد الـ Processors

نعزلها ← Serial

**Workload: sum of 10 scalars, and 10 × 10 matrix sum**

لـ بقدر أجزه

- Speed up from 10 to 100 processors → الـ Processor 1 لنسبة مقارنة

**Single processor: Time = (10 + 100) × $t_{add}$**

رح يحتاج 110 · عملية جمع

**10 processors** م Processor  علية جمع
- Time = 10 × $t_{add}$ + 100/10 × $t_{add}$ = 20 × $t_{add}$ → بدل ما يكون 110 بكون 20
- Speedup = 110/20 = 5.5 (55% of potential)

speedup → S
$Eff = \dfrac{S}{P} * 100\%$

**100 processors** م Processor  علية جمع
- Time = 10 × $t_{add}$ + 100/100 × $t_{add}$ = 11 × $t_{add}$
- Speedup = 110/11 = 10 (10% of potential)

Eff ← بظلها سيئ ← Processors عدد

**Assumes load can be balanced across processors**

MK

---

# Scaling Example (cont)

**What if matrix size is 100 × 100?**

**Single processor: Time = (10 + 10000) × $t_{add}$**

**10 processors**
- Time = 10 × $t_{add}$ + 10000/10 × $t_{add}$ = 1010 × $t_{add}$
- Speedup = 10010/1010 = 9.9 (99% of potential) ← Eff

( وهذا جبر )

**100 processors**
- Time = 10 × $t_{add}$ + 10000/100 × $t_{add}$ = 110 × $t_{add}$
- Speedup = 10010/110 = 91 (91% of potential) ← Eff
- **Assuming load balanced**

MK

# Strong vs Weak Scaling

- Strong scaling: problem size fixed
  - As in example
- Weak scaling: problem size proportional t number of processors
  - 10 processors, 10 × 10 matrix
    - Time = $20 \times t_{add}$
  - 100 processors, 32 × 32 matrix
    - Time = $10 \times t_{add} + 1000/100 \times t_{add} = 20 \times t_{add}$
  - Constant performance in this example

لو كبرنا المسألة
لما نستخدم عدد
أكبر من ال
Processors
بصير weak
Scaling

# Contents

# Contents

# Instruction and Data Streams

An alternate classification

Instruction بيجي لـ

I → [P]

↑
D → Data صتطلعه

|  |  | Data Streams | |
|---|---|---|---|
|  |  | Single | Multiple |
| Instruction Streams | Single | **SISD**: → أبسط اشي<br>Intel Pentium 4 | **SIMD**: SSE (vector) منهم<br>instructions of x86 |
|  | Multiple | **MISD**: org Data بيشتغلوا كتير<br>Instruction بأكتر<br>No examples today | **MIMD**: intel core i7 فيها<br>Intel Xeon e5345 |

برنامج واحد وبيحوي iF statement كل core ينفذ بشانه بس

## SPMD: Single Program Multiple Data

- A parallel program on a MIMD computer

- Conditional code for different processors

كل Processor بشتغل بشغلة مختلفة عن الـ Processor الآخر.

# Vector Processors

- Highly pipelined function units
- <u>Stream</u> data from/to vector registers to units
  - Data collected from memory into registers
  - Results stored from registers to memory
- Example: Vector <u>extension</u> to RISC-V إضافات
  - v0 to v31: 32 × 64-element registers, (64-bit elements)
  - Vector instructions
    - fld.v, fsd.v: load/store vector ( Floating Point )
    - fadd.d.v: add vectors of double
    - fadd.d.vs: add scalar to each element of vector of double

ممكن أضرب Vector ← Scalar → مثلاً
- Significantly reduces instruction-fetch bandwidth
  
  لـ بـ 1 instruction بتنفذ (64 element)

---

# Example: DAXPY (Y = a × X + Y)

① Conventional RISC-V code: بأشر جماية الـ Array

```
        fld     f0,a(x3)        // load scalar a
        addi    x5,x19,512      // end of array X
loop:   fld     f1,0(x19)       // load x[i]
        fmul.d  f1,f1,f0        // a * x[i]
        fld     f2,0(x20)       // load y[i]
        fadd.d  f2,f2,f1        // a * x[i] + y[i]
        fsd     f2,0(x20)       // store y[i]
        addi    x19,x19,8       // increment index to x
        addi    x20,x20,8       // increment index to y
        bltu    x19,x5,loop     // repeat if not done
```

64*8 بأشرك النهاية ← لـ بتنفذوا 64 مرة

2 + 8 * 64 = 514 <= instruction

② Vector RISC-V code: بأشرك

```
        fld         f0,a(x3)        // load scalar a
        fld.v       v0,0(x19)       // load vector x
        fmul.d.vs   v0,v0,f0        // vector-scalar multiply
        fld.v       v1,0(x20)       // load vector y
        fadd.d.v    v1,v1,v0        // vector-vector add
        fsd.v       v1,0(x20)       // store vector y
```

for (i=0, i<64, i++)
  y[i]= a* x[i]+y[i];

6 instruction هذا البرنامج أفضل بس بـ

# Vector vs. Scalar

- Vector architectures and compilers
  - Simplify data-parallel programming
  - Explicit statement of absence of loop-carried dependences
    - Reduced checking in hardware
  - Regular access patterns benefit from interleaved and burst memory
  - Avoid control hazards by avoiding loops
- More general than ad-hoc media extensions (such as MMX, SSE)
  - Better match with compiler technology

# SIMD (Single Instruction Multiple Data)

- Operate elementwise on vectors of data
  - E.g., MMX and SSE instructions in x86
    - Multiple data elements in 128-bit wide registers ) طول يويسم
- All processors execute the same instruction at the same time

  AVx 256 bit
  AVx 512 bit
  ↳
  بنفبر نقسمه
  - Each with different data address, etc.  double ← 8 DP ل
  Single ← 16 SP
- Simplifies synchronization
- Reduced instruction control hardware
- Works best for highly data-parallel applications

Scanned with CamScanner

# Vector vs. Multimedia Extensions

متثل SIMD

- Vector instructions have a variable vector width, multimedia extensions have a fixed width

- Vector instructions support strided access, multimedia extensions do not

- Vector units can be combination of pipelined and arrayed functional units:



Lane 0  Lane 2  Lane 3  Lane 4

FU Pipelined

Element group

*(handwritten notes in Arabic surrounding the figure)*

4 FU بدل ما كل ... FU

register التقسيم الـ
لـ 4 أجزاء

# Contents

بسهل ونقلل عدد الـ input/output الـ register الواحد

# Multithreading

- Performing multiple threads of execution in parallel → Multiple Process   واحد ينفذ Core
  - Replicate registers, PC, etc.
  - Fast switching between threads
- Fine-grain multithreading
  - Switch threads after each cycle
  - Interleave instruction execution
  - If one thread stalls, others are executed
- Coarse-grain multithreading
  - Only switch on long stall (e.g., L2-cache miss)
  - Simplifies hardware, but doesn't hide short stalls (eg, data hazards)

# Simultaneous Multithreading

- In multiple-issue dynamically scheduled processor (SMT)
  - Schedule instructions from multiple threads
  - Instructions from independent threads execute when function units are available
  - Within threads, dependencies handled by scheduling and register renaming
- Example: Intel Pentium-4 HT → Hyper Threading
  - Two threads: duplicated registers, shared function units and caches

# Multithreading Example



Issue slots ⟶

Thread A   Thread B   Thread C   Thread D

Time

* يبين إنه الـ Thread بحالات نادرة جدًا بغيروا ينفذوا الـ 4

Issue slots ⟶

① Coarse MT   ② Fine MT   ③ SMT

Time

لما يشعر أنه في توقف بينتقل لله ثم الـ Thread اللي بعده، قادر على إخفاء الـ stalls وهما

بختلطوا الـ Threads مع بعض بس كل وحدة موجا وهي الـ Flag هي لأي الـ Thread نابعة وبختلط تنفيذهم بس مشي على شكل مجموعات. أعقد

له كل cycle بجيب من Thread

وهاد كويس لإنه ببعد الـ Instruction بـ thread واحد عن بعض وهاد بحسن لو كان في dependence، قادر على إخفاء Short /small Stalls أبسط

**MK**

# Future of Multithreading

- Will it survive? In what form?
- Power considerations ⇒ simplified microarchitectures
  - Simpler forms of multithreading
- Tolerating cache-miss latency
  - Thread switch may be most effective
- Multiple simple cores might share resources more effectively

# Contents

# Shared Memory

SMP: shared memory multiprocessor

- Hardware provides single physical address space for all processors
- Synchronize shared variables using locks
- Memory access time
  - ① UMA (uniform) vs. NUMA (nonuniform) → Shared أنواع الـ Memory

*لا بجيبين نفس البعد*    *كل Mem الـ* ②



→ 1 Memory نشبكهم مع

§6.5 Multicore and Other Shared Memory Multiprocessors

# Example: Sum Reduction

Sum 64,000 numbers on 64 processor UMA

- Each processor has ID: $0 \leq Pn \leq 63$    كشان ان
  
     64 processor
- Partition 1000 numbers per processor
- Initial summation on each processor

```
sum[Pn] = 0;
  for (i = 1000*Pn;
       i < 1000*(Pn+1); i += 1)
    sum[Pn] += A[i];
```

Now need to add these partial sums

- Reduction: divide and conquer
- Half the processors add pairs, then quarter, …
- Need to synchronize between reduction steps

# Example: Sum Reduction

( التصنيف ودكون مسؤولية النص التحتاني ،انه بينقل قيمة الفوقاني )



(half = 1) [0] [1]

(half = 2) [0] [1] [2] [3]

(half = 4) [0] [1] [2] [3] [4] [5] [6] [7]

```
half = 64;
do
  synch();
  if (half%2 != 0 && Pn == 0)
    sum[0] += sum[half-1];
    /* Conditional sum needed when half is odd;
       Processor0 gets missing element */
  half = half/2; /* dividing line on who sums */
  if (Pn < half) sum[Pn] += sum[Pn+half];
while (half > 1);
```

الهدف ← انا كان الرقم ودي كشان نعالج المشكلة

# Contents

# History of GPUs

- Early video cards بالمقاييات
  - Frame buffer memory with address generation for video output
- 3D graphics processing بنهاية الثماينات
  - Originally high-end computers (e.g., SGI)
  - Moore's Law ⇒ lower cost, higher density
  - 3D graphics cards for PCs and game consoles مثال البلايستيشن
- Graphics Processing Units GPUs
  - Processors oriented to 3D graphics tasks
  - Vertex/pixel processing, shading, texture mapping, rasterization تحدد الأماكن المحتانات

تحويل الفيديو لـ Signals

# Graphics in the System ( Computer )



اقتلنا لهذا التصميم

انتقلنا لهذا تكلفته أقل

*Intel CPU* — x16 PCI-Express Link — Front Side Bus — GPU — North Bridge — DDR2 Memory — display — x4 PCI-Express Link derivative — 128-bit 667 MT/s — GPU Memory — South Bridge

CPU — Front Side Bus — North Bridge — Memory — PCI Bus — South Bridge — VGA Controller — Framebuffer Memory — LAN — UART — VGA Display

AMD CPU — CPU core — Internal bus — North Bridge — 128-bit 667 MT/s — DDR2 Memory — x16 PCI-Express Link — HyperTransport 1.03 — GPU — Chipset — display — GPU Memory

"تحت قديمة"

# GPU Architectures

→ Alot of data

- Processing is highly data-parallel لغنى الحسابات
  - GPUs are highly multithreaded
  - Use thread switching to hide memory latency
    - Less reliance on multi-level caches
  - Graphics memory is wide and high-bandwidth
- Trend toward general purpose GPUs   متنوعة التطبيقات
  - Heterogeneous CPU/GPU systems   → Performance لتعطيني أفضل من الـ CPU
  - CPU for sequential code, GPU for parallel code
- Programming languages/APIs
  - DirectX, OpenGL
  - C for Graphics (Cg), High Level Shader Language (HLSL)
  - Compute Unified Device Architecture (CUDA)
    نوعه من NVidia

# Modern Computer

*i 3*        *i 7*

*Nvidia*

## INTEL® Z390 CHIPSET BLOCK DIAGRAM

**1x16 lanes PCI Express* 3.0 Graphics or Intel SSD**

OR

**2x8 lanes PCI Express* 3.0 Graphics and Intel SSD**

OR

**1x8 and 2x4 lanes PCI Express* 3.0 Graphics and Intel SSD**

**8th Gen Intel® Core™ Processors**

Intel® UHD Graphics

**DDR4 2xDIMMs per Channel Up to 2666 MHz[1]**

**DDR4 2xDIMMs per Channel Up to 2666 MHz[1]**

**Three Independent DP/HDMI Display Support**

DMI 3.0

**Intel® Optane™ Memory Support[3]**

**Up to 24 x PCI Express* 3.0**

8 Gb/s each x 1

**Intel® Smart Sound Technology[1]**

**6 x SATA 6 Gb/s Ports; SATA Port Disable**

Up to

Intel Z390

**Intel® High Definition Audio[1]**

Intel® Z390 Chipset block diagram

- PCI Express* 3.0 Graphics, and Intel SSD
- Three Independent DP/HDMI Display Support
- Up to 24 x PCI Express* 3.0 — 8 Gb/s each x 1
- 6 x SATA 6 Gb/s Ports; SATA Port Disable — Up to 6 Gb/s
- Up to 6 x USB 3.1 Gen 2 Ports; Up to 10 x USB 3.1 Gen 1 Ports; 14 x USB 2.0 Ports
- Intel® Integrated 10/100/1000 MAC
  - PCIe* x 1
  - SMBus
  - Intel® Ethernet Connection

DMI 3.0 — Intel® Z390 Chipset — SPI

- Intel® Optane™ Memory Support¹
- Intel® Smart Sound Technology¹
- Intel® High Definition Audio¹
- Intel® Rapid Storage Technology with RAID¹
- Intel® Rapid Storage Technology for PCI Express* Storage¹
- Intel® Wireless-AC 802.11ac and Bluetooth* 5
- Intel® Wireless-AC Adapter

- Intel® ME Firmware
- Intel® Platform Trust Technology¹
- Intel® Extreme Tuning Utility Support

Optional

# Example: NVIDIA Fermi

Multiple SIMD processors, each as shown:



كل وحدة بتنفذ ← Instruction نفسها لكن على Data مختلفة

SIMD Lanes (Thread Processors)

Address coalescing unit    Interconnection network

Local Memory 64 KiB    ( داخل GPU )

To Global Memory (خارج GPU)

# Example: NVIDIA Fermi

- SIMD Processor: 16 SIMD lanes
- SIMD instruction
  - Operates on 32 element wide threads
  - Dynamically scheduled on 16-wide processor over 2 cycles
- 32K x 32-bit registers spread across lanes
  - 64 registers per thread context

# GPU Memory Structures

CUDA Thread

Per-CUDA Thread Private Memory

Thread block

Per-Block Local Memory

Grid 0

Grid 1

— — — Inter-Grid Synchronization — — —

GPU Memory

Sequence

*(handwritten Arabic notes:)*

مسؤولية الـ cpu قبل ما نحطي المجال

الـ cpu ، إذا تعمل الحسابات المطلوبة

إذا الـ data موجودة بالـ Main Mem لازم

تنتقل الـ cpu بعدين بيتم عليها الحسابات بسرعة

ثم تنتقل النتائج من

أخرى الـ Mem

CPU → M

---

# Classifying GPUs

- Don't fit nicely into SIMD/MIMD model
  - Conditional execution in a thread allows an illusion of MIMD
    - But with performance degredation
    - Need to write general purpose code with care

|  | Static: Discovered at Compile Time | Dynamic: Discovered at Runtime |
|---|---|---|
| Instruction-Level Parallelism | VLIW | Superscalar |
| Data-Level Parallelism | SIMD or Vector | **Tesla Multiprocessor** GPUs |

# INTRODUCING TURING

## TU102 – FULL CONFIG
18.6 BILLION TRANSISTORS

| | |
|---|---|
| SM | 72 |
| CUDA CORES | 4608 |
| TENSOR CORES | 576 |
| RT CORES | 72 |
| GEOMETRY UNITS | 36 |
| TEXTURE UNITS | 288 |
| ROP UNITS | 96 |
| MEMORY | 384-bit 7 GHz GDDR6 |
| NVLINK CHANNELS | 2 |

Shapes    Ink to Shape

# RTX 2080 Ti

Thursday, April 23, 2020      9:53 AM

# INTRODUCING TURING

## TU102 – FULL CONFIG

### 18.6 BILLION TRANSISTORS

| SM | 72 |
|---|---|
| CUDA CORES | 4608 |
| TENSOR CORES | 576 |
| RT CORES | 72 |
| GEOMETRY UNITS | 36 |

# Putting GPUs into Perspective

| Feature | Multicore with SIMD *(CPU)* | GPU |
|---|---|---|
| SIMD processors | 4 to 8 | 8 to 16 |
| SIMD lanes/processor | 2 to 4 | 8 to 16 *or 32 or 64* |
| Multithreading hardware support for SIMD threads | 2 to 4 | 16 to 32 |
| Typical ratio of single precision to double-precision performance | 2:1 | 2:1 |
| Largest cache size | 8 MB | 0.75 MB |
| Size of memory address | 64-bit | 64-bit |
| Size of main memory | 8 GB to 256 GB | 4 GB to 6 GB |
| Memory protection at level of page | Yes | Yes |
| Demand paging | Yes | No |
| Integrated scalar processor/SIMD processor | Yes | No |
| Cache coherent | Yes | No |

# Guide to GPU Terms

| Type | More descriptive name | Closest old term outside of GPUs | Official CUDA/NVIDIA GPU term | Book definition |
|---|---|---|---|---|
| Program abstractions | Vectorizable Loop | Vectorizable Loop | Grid | A vectorizable loop, executed on the GPU, made up of one or more Thread Blocks (bodies of vectorized loop) that can execute in parallel. |
| | Body of Vectorized Loop | Body of a (Strip-Mined) Vectorized Loop | Thread Block | A vectorized loop executed on a multithreaded SIMD Processor, made up of one or more threads of SIMD instructions. They can communicate via Local Memory. |
| | Sequence of SIMD Lane Operations | One iteration of a Scalar Loop | CUDA Thread | A vertical cut of a thread of SIMD instructions corresponding to one element executed by one SIMD Lane. Result is stored depending on mask and predicate register. |
| Machine object | A Thread of SIMD Instructions | Thread of Vector Instructions | Warp | A traditional thread, but it contains just SIMD Instructions that are executed on a multithreaded SIMD Processor. Results stored depending on a per-element mask. |
| | SIMD Instruction | Vector Instruction | PTX Instruction | A single SIMD instruction executed across SIMD Lanes. |
| Processing hardware | Multithreaded SIMD Processor | (Multithreaded) Vector Processor | Streaming Multiprocessor | A multithreaded SIMD Processor executes threads of SIMD instructions, independent of other SIMD Processors. |
| | Thread Block Scheduler | Scalar Processor | Giga Thread Engine | Assigns multiple Thread Blocks (bodies of vectorized loop) to multithreaded SIMD Processors. |
| | SIMD Thread Scheduler | Thread scheduler in a Multithreaded CPU | Warp Scheduler | Hardware unit that schedules and issues threads of SIMD instructions when they are ready to execute; includes a scoreboard to track SIMD Thread execution. |
| | SIMD Lane | Vector lane | Thread Processor | A SIMD Lane executes the operations in a thread of SIMD instructions on a single element. Results stored depending on mask. |
| Memory hardware | GPU Memory | Main Memory | Global Memory | DRAM memory accessible by all multithreaded SIMD Processors in a GPU. |
| | Local Memory | Local Memory | Shared Memory | Fast local SRAM for one multithreaded SIMD Processor, unavailable to other SIMD Processors. |
| | SIMD Lane Registers | Vector Lane Registers | Thread Processor Registers | Registers in a single SIMD Lane allocated across a full thread block (body of vectorized loop). |

# Contents

# Message Passing

- Each processor has private physical address space

- Hardware sends/receives messages between processors



| Processor | Processor | . . . | Processor |
|---|---|---|---|
| Cache | Cache | . . . | Cache |
| Memory | Memory | . . . | Memory |

Interconnection Network

# Loosely Coupled Clusters

ما بتعاونوا عن طريق الـ Memory

- Network of independent computers
  - Each has private memory and OS
  - Connected using I/O system
    - E.g., Ethernet/switch, Internet → بطيئ
- Suitable for applications with independent tasks
  - Web servers, databases, simulations, ...
- High availability, scalable, affordable
- Problems ↳  انا واحد تعطل البقية بكملوا . بستخدموه معوا عشان يسول
  - Administration cost (prefer virtual machines)
  - Low interconnect bandwidth ككس
    - c.f. processor/memory bandwidth on an SMP → هدول أفضل

# Sum Reduction (Again)

- Sum 64,000 on 64 processors → 64000 موزعة  64 P كل P بياخذ 1000
- First distribute 1000 numbers to each
  - The do partial sums

```
sum = 0;
for (i = 0; i<1000; i += 1)
    sum += AN[i];
```

- Reduction
  - Half the processors send, other half receive and add
  - The quarter send, quarter receive and add, ...

# Sum Reduction (Again)

- Given send() and receive() operations

```
limit = 64; half = 64;/* 64 processors */
do
   half = (half+1)/2; /* send vs. receive
                         dividing line */
   if (Pn >= half && Pn < limit)
     send(Pn - half, sum);
   if (Pn < (limit/2))
     sum += receive();
   limit = half; /* upper limit of senders */
while (half > 1); /* exit with final sum */
```

هذا الجديد المهم

- Send/receive also provide synchronization
- Assumes send/receive take similar time to addition

# Grid Computing

- Separate computers interconnected by long-haul networks
  - E.g., Internet connections
  - Work units farmed out, results sent back
- Can make use of idle time on PCs
  - E.g., SETI@home, World Community Grid

* يعني لو الجهاز هو بستخدمه دقتر انتجاه عالانترنت لحد ثم بيستخدمه *

Scanned with CamScanner

# Contents

نشيك الـ Processors أو الـ computerss مع بعض

# Interconnection Networks

٭خاصية الـ 2D Torus عن الـ Mesh اللي بالصورة

.انه يسرع علية التواصل لإنه نشبك الأطراف مع بعض

و قلل الـ latency .

- Network topologies
  - Arrangements of processors, switches, and links

① 

Bus

② 

Ring

③ 

2D Mesh
(Torus)

④ 

N-cube (N = 3)

diminsion = log N

latency = diminsion

⑤ 

$\frac{N(N-1)}{2}$

مشكلته مكلفة جدا وعدد الأسلاك كبير

Fully connected  latency = 1

كل نقطة مشبوكة مع الكل

# Interconnection Networks

Mesh Topology

2D Torus

3D Torus

1D

2D

3D

4D

# Multistage Networks



a. Crossbar

" مكلفة جداً الـ Switch غالي "

b. Omega network → لحتى يقللوا من التكلفة

الـ Switch box لبستخدموا

c. Omega network switch box

# Network Characteristics

Performance
- Latency per message (unloaded network) بيهمش بحساب الـ cost (أقل مايمكن)
- Throughput
  - Link bandwidth → Frequency والـ الأسلاك وعدد على بعتمد (بنفرح اكتر ما يمكن بس مش محسوب)
  - Total network bandwidth
  - Bisection bandwidth → عندي بظل link كم لنعين أقسمهم لما (أعلى ما يمكن بس مش بحساب الـ cost)
- Congestion delays (depending on traffic)

Cost ( معقولة )

Power

Routability in silicon

Handwritten annotations at top: a box with "8" and "8" labels; $d = \lg 14 = 6$ ; "?" ; $\frac{64 \times 63}{2}$

| Evaluation category | Bus | Ring | 2D mesh | 2D torus | Hypercube | Fat tree | Fully connected |
|---|---|---|---|---|---|---|---|
| **Performance** | | | | | | | |
| BW$_{\text{Bisection}}$ in # links | 1 | 2 | 8 | 16 | 32 | 32 | 1024 |
| Max (ave.) hop count | 1 (1) | 32 (16) | 14 (7) | 8 (4) | 6 (3) | 11 (9) | 1 (1) |
| **Cost** | | | | | | | |
| I/O ports per switch | NA | 3 | 5 | 5 | 7 | 4 | 64 |
| Number of switches | NA | 64 | 64 | 64 | 64 | 192 | 64 |
| Number of net. links | 1 | 64 | 112 | 128 | 192 | 320 | 2016 |
| Total number of links | 1 | 128 | 176 | 192 | 256 | 384 | 2080 |

**Figure F.15** Performance and cost of several network topologies for 64 nodes. The bus is the standard reference at unit network link cost and bisection bandwidth. Values are given in terms of bidirectional links and ports. Hop count includes a switch and its output link, but not the injection link at end nodes. Except for the bus, values are given for the number of network links and total number of links, including injection/reception links between end node devices and the network.

CH6 :

Slide 9 : Amdahl's law :

$T\_old = T\_old * (1-f+f)$

$T\_new = T\_old * (1-f + f/p)$     → عدد الProcessors

$Speed\ up = T\_old / T\_new = 1/(1-f+f/p)$

Serial →    Told    → Parallelizable

← زمن قديم وكبير    | 1-f | f |

Parallel    ← بدي أوصل لزمن جديد عن طريق الـ    | 1-f | f/p |
Program                       Tnew

* Serial :    لأبد إنك تنفذه هما بتقدر تعدله

* parallelizable :    ممكن انفذه لأنه ينفذه وأكثر من P

∴ $\lim\limits_{P\to\infty} Speed\ up = \dfrac{1}{1-f}$    → أعلى Speedup ممكن تحصل عليها

Slide 11 : Strong Scalling    ( تخلي المسألة ثابتة وتزيد عدد الـ Processors )

| Size | | | P=1 | | P=10 | | P=100 | |
|---|---|---|---|---|---|---|---|---|
| 10+100 | Speedup | Efficiency | 1 | 100% | 5.5 | 55% | 10 | 10% |
| 10+1000 | Speedup | Efficiency | 1 | 100% | 9.2 | 92% | 51 | 51% |
| 10+10000 | Speedup | Efficiency | 1 | 100% | 9.9 | 99% | 91 | 91% |

إضافة من الدكتور →

* بنلاحظ كل ما زاد عدد الـ Processors الـ Efficiency بتقل لأن الـ المسألة ثابتة.

Slide 16 : Vector processor ( رسمة توضيحية )

Main Memory

Vector(load/ Store)

FP add/subtract

FP multiply

FP divide

Integer

Logical

Functional unit

Vector $^{64}$ registers

Scalar registers

* المهم هنا, انه داخل الـ Vector computer في Vector registers وهذه الـ Vectors ممكن يكون فيه منزم 32 لكن كل واحد منزم واحد

Vector register مثل بشبح لـ value 1 لـ بشبح لـ Vector كامل فكل واحد من الـ Vector register ممكن يكون بتكون من ( 64 register )

⇒ 32 VR × 64 = 2K reg

( مثلاً هياكل هو غالي )

* Scalar ممكن يكون 32 واحد بس كل واحد

Registe : فيه value 1

## Slide 19 : SIMD          ( add VC, VA, VB )

| A.0 | A.1 | A.2 | A.3 | VA |
|-----|-----|-----|-----|-----|
| B.0 | B.1 | B.2 | B.3 | VB |
| + | + | + | + | |
| C.0 | C.1 | C.2 | C.3 | VC |

\* يحسن الـ Performence بحيث ، 1 instruction بيعمل 4 additions .

## Slide 22 : Multithreading

PC → I\$ → RF → FUs

PC₂ → RF₂ →

## Slide 27 : NUMA

| Memory | CPU | CPU | | CPU | CPU | Memory |
|--------|-----|-----|---|-----|-----|--------|
| | CPU | CPU | > | CPU | CPU | |

Bus Interconnect

| Memory | CPU | CPU | | CPU | CPU | Memory |
|--------|-----|-----|---|-----|-----|--------|
| | CPU | CPU | < | CPU | CPU | |

\* كل مجموعة ، الها جزء خاص من الـ Memory و الـ 4 Memories منتشاركين وأي CPU
بتشوف الباقي لكن الـ Memory اللي في الـ 4 cpus قريبة أكثر من باقي الـ Memory وبتوصلولها
بشكل أسرع .

# Slide 28 : Sum Reduction

## Serial :

```
Sum = 0;
for (i=0; i<64000; i += 1)
    Sum += A[i];
```

## Parallel :

```
Sum [Pn] = 0;
for (i=1000 * Pn; i < 1000* (Pn+1); i+=1)
    Sum [Pn] += A[i];
```

عشان نخذمر الوقت، بنقسم كل الشغل ←
على (64 Processor). حنخلص أقل بكثير من الـ Serial
بطلعلي 64 جواب و بعد بجمعهم كلهم برقم واحد
وبدل ما أعمل loop عشان أجمع الـ 64 بعمل الي بسلايد 29

---

# Slide 41 : Message-passing Multicomputers

# Parallel Benchmarks

- *Linpack:* matrix linear algebra
- *SPECrate:* parallel run of SPEC CPU programs
  - Job-level parallelism
- SPLASH: Stanford Parallel Applications for Shared Memory
  - Mix of kernels and applications, strong scaling
- NAS (NASA Advanced Supercomputing) suite
  - computational fluid dynamics kernels
- PARSEC (Princeton Application Repository for Shared Memory Computers) suite
  - Multithreaded applications using Pthreads and OpenMP

# Code or Applications?

- Traditional benchmarks
  - Fixed code and data sets
- Parallel programming is evolving
  - Should algorithms, programming languages, and tools be part of the system?
  - Compare systems, provided they implement a given application
  - E.g., Linpack, Berkeley Design Patterns
- Would foster innovation in approaches to parallelism

بقوي الابتكار بالطرق
اللي الها علاقة بال
Parallelism

و الهدف نقيم الـperformance ونعرخاهل هو كويس ولامحتاج الكتب برنامج أحسنه

# Modeling Performance

- Assume performance metric of interest is achievable GFLOPs/sec

قديوه computer بقدر يعمل
وعش داياهو المعيار الوحيد بين هو الأهم هون   Floating Point operations

  - Measured using computational kernels from Berkeley Design Patterns
- Arithmetic intensity of a kernel

موضوع
بين الـKernel
المختلفات

  - FLOPs per byte of memory accessed
- For a given computer, determine

ALU s الـكل لو شغلنا لو هـ
Computer الـ هذا يعني.
قديش ممكن نعمل
GFLOPs

  - Peak GFLOPS (from data sheet)
  - Peak memory bytes/sec (using Stream benchmark)

# Roofline Diagram



*logarithmic scale*

Attainable GFLOPs/second axis: 0.5, 1.0, 2.0, 4.0, 8.0, 16.0, 32.0, 64.0

peak memory BW (stream) / peak floating-point performance

Kernel 1 (Memory Bandwidth limited)

Kernel 2 (Computation limited)

Arithmetic Intensity: FLOPs/Byte Ratio — 1/8, 1/4, 1/2, 1, 2, 4, 8, 16

Attainable GPLOPs/sec
= Max ( Peak Memory BW × Arithmetic Intensity, Peak FP Performance )

# Comparing Systems

- Example: Opteron X2 vs. Opteron X4
  - 2-core vs. 4-core, 2× FP performance/core, 2.2GHz vs. 2.3GHz, 1 × 2 SIMD vs. 2 × 2 SIMD

  بعض 2 SIMD op كل cycle لـ        لـ كل cycle بنفذ منها 1 SIMD بعمليتين
  كل وحدة فيها عمليتين

  - Same memory system



Opteron X4 (Barcelona)

Opteron X2

ثابت للـين

Attainable GFLOP/s axis: 0.5, 1.0, 2.0, 4.0, 8.0, 16.0, 32.0, 64.0, 128.0

Actual FLOPbyte ratio — 1/8, 1/4, 1/2, 1, 2, 4, 8, 16

To get higher performance on X4 than X2
- Need high arithmetic intensity
- Or working set must fit in X4's 2MB L-3 cache

opteron X2 : Peak FP Performance = 2 cores/chip * 2 SIMD/core *2 FP/SIMD * 2.2 Gcycle/s = 2*2 *2*2.2 = 17.6 GFlops /s

opteron X4 : Peak FP Performance = 4 cores/chip * 4 SIMD/core * 2 FP /SIMD * 2-3 Gcycle/s = 4*4*2 * 2.3 = 73.6 GFlops /s

# Optimizing Performance

- Optimize FP performance
  - Balance adds & multiplies
  - Improve superscalar ILP and use of SIMD instructions *(Instruction level Parallelism)*
- Optimize memory usage
  - Software prefetch
    - Avoid load stalls
  - Memory affinity
    - Avoid non-local data accesses



AMD Opteron



AMD Opteron

---

# Optimizing Performance

- Choice of optimization depends on arithmetic intensity of code



Arithmetic Intensity: FLOPs/Byte Ratio

- Arithmetic intensity is not always fixed
  - May scale with problem size
  - Caching reduces memory accesses
    - Increases arithmetic intensity

Scanned with CamScanner

# Contents

# i7-960 vs. NVIDIA Tesla 280/480

|  | CPU | GPU | GPU |  |  |
|---|---|---|---|---|---|
|  | Core i7-960 | GTX 280 | GTX 480 | Ratio 280/17 | Ratio 480/17 |
| Number of processing elements (cores or SMs) | 4 | 30 | 15 | 7.5 | 3.8 |
| Clock frequency (GHz) | 3.2 | 1.3 | 1.4 | 0.41 | 0.44 |
| Die size | 263 | 576 | 520 | 2.2 | 2.0 |
| Technology | Intel 45 nm | TCMS 65 nm | TCMS 40 nm | 1.6 | 1.0 |
| Power (chip, not module) | 130 | 130 | 167 | 1.0 | 1.3 |
| Transistors | 700 M | 1400 M | 3100 M | 2.0 | 4.4 |
| Memory brandwith (GBytes/sec) | 32 | أكبر 141 | 177 | 4.4 | 5.5 |
| Single frecision SIMD width | 4 | 8 | 32 | 2.0 | 8.0 |
| Dobule precision SIMD with | أحسن 2 | أسوء 1 | 16 | 0.5 | 8.0 |
| Peak Single frecision scalar FLOPS (GFLOP/sec) | 26 | 117 | 63 | 4.6 | 2.5 |
| Peak Single frecision s SIMD FLOPS (GFLOP/Sec) | 102 | 311 to 933 | 515 to 1344 | 3.0-9.1 | 6.6-13.1 |
| (SP 1 add or multiply) | N.A. | (311) | (515) | (3.0) | (6.6) |
| (SP 1 instruction fused) | N.A | (622) | (1344) | (6.1) | (13.1) |
| (face SP dual issue fused) | N.A | (933) | N.A | (9.1) | – |
| Peal double frecision SIMD FLOPS (GFLOP/sec) | 51 | 78 | 515 | 1.5 | 10.1 |

# Rooflines

high performence اسهل تشغيل

# Benchmarks

| Kernel | Units | Core I7-960 | GTX 280 | GTX 280/ I7-960 |
|--------|-------|-------------|---------|-----------------|
| SGEMM | GFLOP/sec | 94 | 364 | 3.9 |
| MC | Billion paths/sec | 0.8 | 1.4 | 1.8 |
| Conv | Million pixels/sec | 1250 | 3500 | 2.8 |
| FFT | GFLOP/sec | 71.4 | 213 | 3.0 |
| SAXPY | GBytes/sec | 16.8 | 88.8 | 5.3 |
| LBM | Million lookups/sec | 85 | 426 | 5.0 |
| Solv | Frames/sec | 103 | 52 | 0.5 |
| SpMV | GFLOP/sec | 4.9 | 9.1 | 1.9 |
| GJK | Frames/sec | 67 | 1020 | 15.2 |
| Sort | Million elements/sec | 250 | 198 | 0.8 |
| RC | Frames/sec | 5 | 8.1 | 1.6 |
| Search | Million queries/sec | 50 | 90 | 1.8 |
| Hist | Million pixels/sec | 1517 | 2583 | 1.7 |
| Bilat | Million pixels/sec | 83 | 475 | 5.7 |

Scanned with CamScanner

# Performance Summary

GPU (480) has 4.4 X the memory bandwidth
- Benefits memory bound kernels

GPU has 13.1 X the single precision throughout, 2.5 X the double precision throughput
- Benefits FP compute bound kernels

CPU cache prevents some kernels from becoming memory bound when they otherwise would on GPU

GPUs offer scatter-gather, which assists with kernels with strided data

Lack of synchronization and memory consistency support on GPU limits performance for some kernels

# Contents

Scanned with CamScanner

# Fallacies

- Amdahl's Law doesn't apply to parallel computers
  - Since we can achieve linear speedup
  - But only on applications with weak scaling
- Peak performance tracks observed performance
  - Marketers like this approach!
  - But compare Xeon with others in example
  - Need to be aware of bottlenecks

# Pitfalls

- Not developing the software to take account of a multiprocessor architecture
  - Example: using a single lock for a shared composite resource
    - Serializes accesses, even if they could be done in parallel
    - Use finer-granularity locking
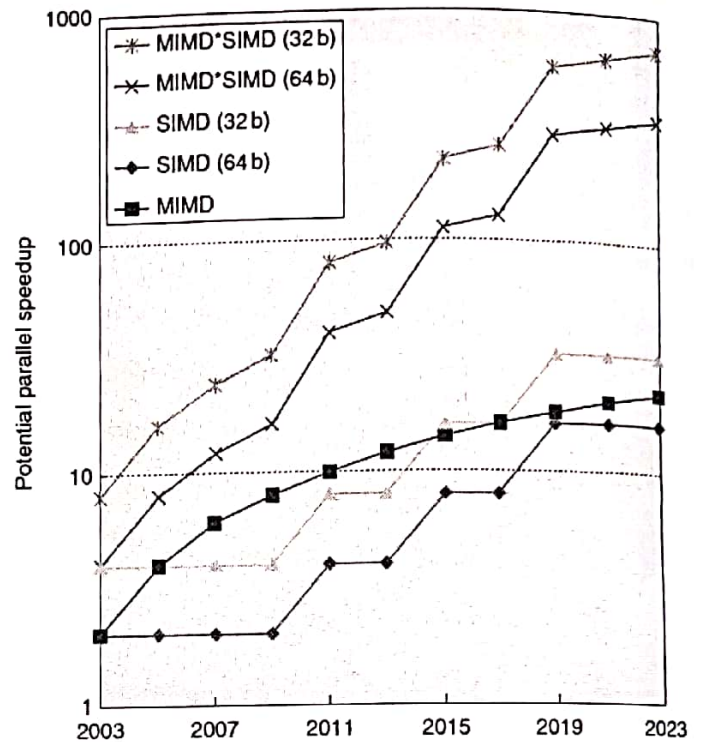
# Contents

# Concluding Remarks

- Goal: higher performance by using multiple processors
- Difficulties
  - Developing parallel software
  - Devising appropriate architectures
- SaaS importance is growing and clusters are a good match
- Performance per dollar and performance per Joule drive both mobile and WSC

# Concluding Remarks (con't)

- SIMD and vector operations match multimedia applications and are easy to program

- Adding 2 cores/chip every 2 years.
- Doubling SIMD operations every 4 years.

# Slide 53 : Roofline Model

## DAXPY :

```
for (i=0; i<N; i+=1)
    Y[i] += S * X[i] + Y[i];
```

∴ Arithmetic Intensity = 2 Flops / (2 loads + 1 store)

$= 2 / (3*8) = 0.0833$ Flop/byte

8 byte كلواحد

( Memory intensive قليل فواحنا )