# Master Informatics Eng.

2018/19

*A.J.Proença*

## Concepts from undegrad Computer Systems

*(most slides are borrowed, mod's in green)*
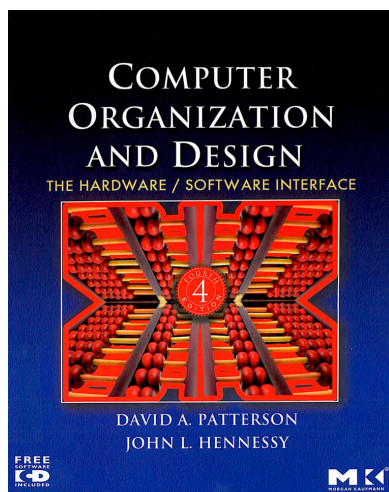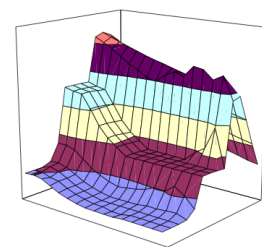
---

## Concepts from undegrad Computer Systems

*– most slides are borrowed from*

*and some from*

Computer Systems
*A Programmer's Perspective* [1]
(*Beta Draft*)

COMPUTER
ORGANIZATION
AND DESIGN
THE HARDWARE / SOFTWARE INTERFACE

4

DAVID A. PATTERSON
JOHN L. HENNESSY

MK

Randal E. Bryant
David R. O'Hallaron

August 1, 2001

*more details at*
*http://gec.di.uminho.pt/miei/sc/*
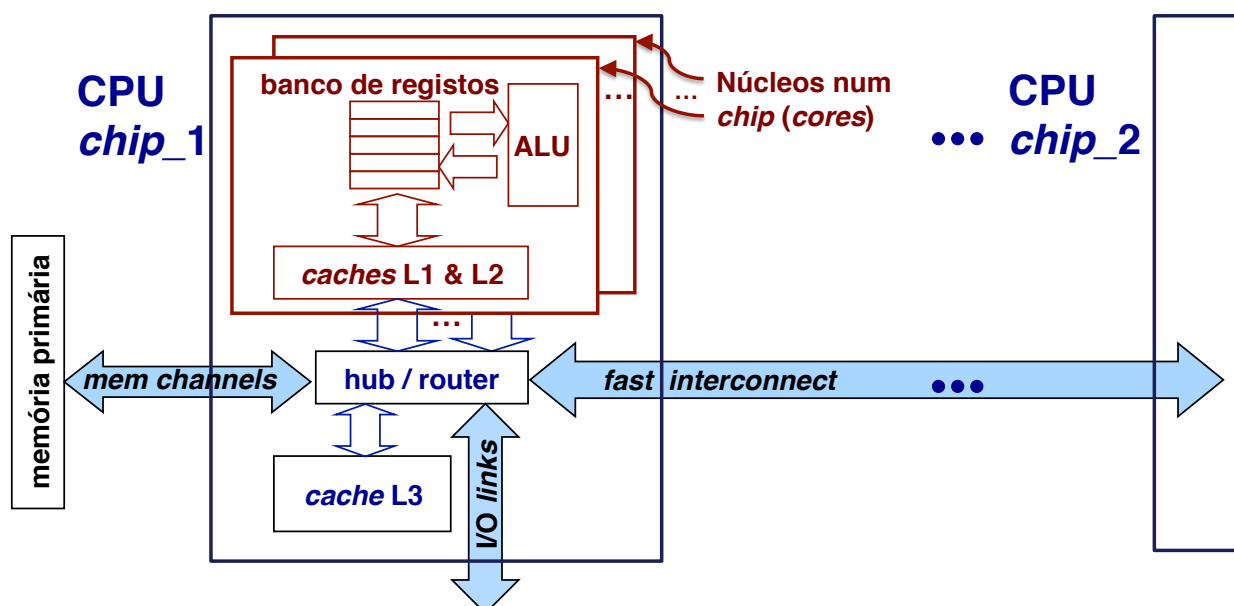
# Background for Advanced Architectures

## Key concepts to revise:

– *numerical data representation (for error analysis)*

– *ISA (Instruction Set Architecture)*

– *how C compilers generate code (a look into assembly code)*
  - *how scalar and structured data are allocated*
  - *how control structures are implemented*
  - *how to call/return from function/procedures*
  - *what architecture features impact performance*

– *Improvements to enhance performance in a <u>single CPU</u>*
  - *ILP: pipeline, multiple issue, …*
  - *data parallelism: SIMD/vector processing, ...*
  - *thread-level parallelism*
  - *memory hierarchy: cache levels, ...*
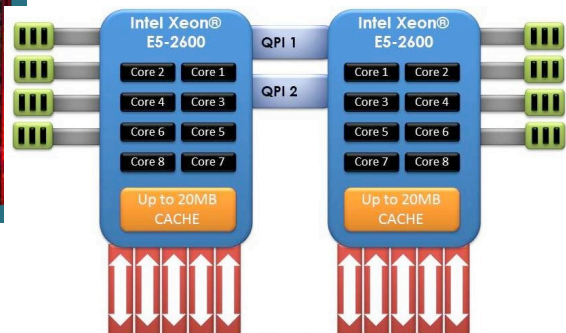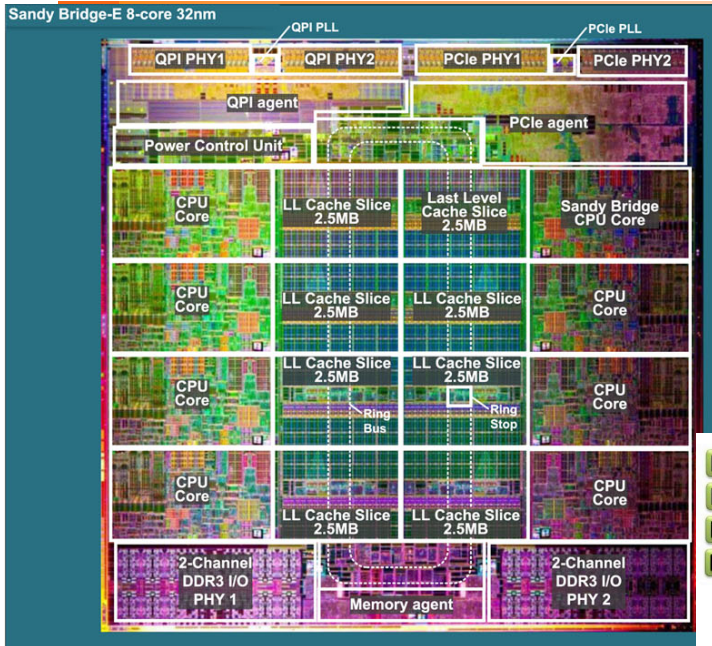
# A hierarquia de cache em arquiteturas multicore

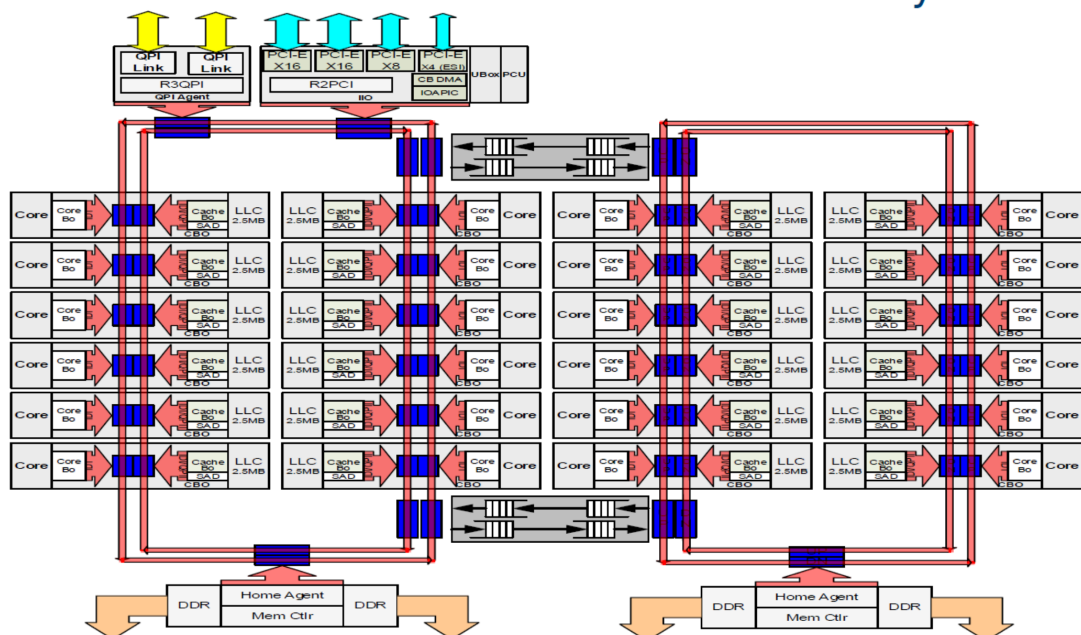## As arquiteturas *multicore* mais recentes:

*AJProença, Sistemas de Computação, UMinho, 2013/14*

**Intel in 2016:**
**Broadwell-EP Xeon** *(22-core)*

## Intel® Xeon® Processor E5 v4 Product Family HCC

**Chips da Intel em 2012/13:**
**Xeon Phi com 60 cores**

PCIe I/O

PCIe I/O Logic

| Core | Core | ... | Core | Core |
| L2 Cache | L2 Cache | | L2 Cache | L2 Cache |

Bidirectional Ring Bus

GDDR5 Memory Controllers

Bidirectional Ring Bus

GDDR5 Memory Controllers

MEMORY I/O

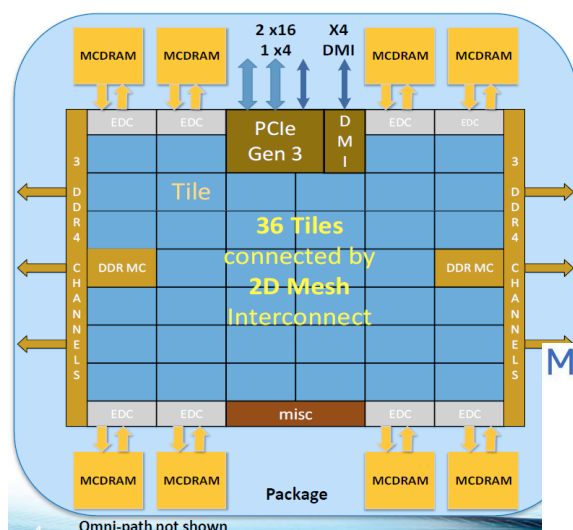| Core | Core | ... | Core | Core |
| L2 Cache | L2 Cache | | L2 Cache | L2 Cache |

*AJProença, Sistemas de Computação, UMinho, 2013/14*      7

---

**Intel new Phi in 2016:**
**KNL with 72 cores**

# Knights Landing Overview

**TILE**

| 2 VPU | CHA | 2 VPU |
| | 1MB L2 | |
| Core | | Core |



2 x16
1 x4

X4 DMI

PCIe Gen 3

Tile

**36 Tiles**
connected by
**2D Mesh**
Interconnect

DDR MC    DDR MC

**Chip: 36 Tiles** interconnected by **2D Mesh**

**Tile**: 2 Cores + 2 VPU/core + 1 MB L2

**Memory: MCDRAM:** 16 GB on-package; High BW

         **DDR4:** 6 channels @ 2400  up to 384GB
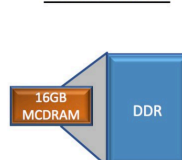
**IO:** 36 lanes PCIe Gen3. 4 lanes of DMI for chipset

**Node:** 1-Socket only
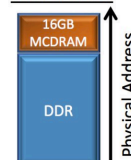
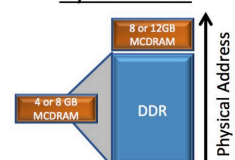**Fabric:** Omni-Path on-package (not shown)

## Memory Modes

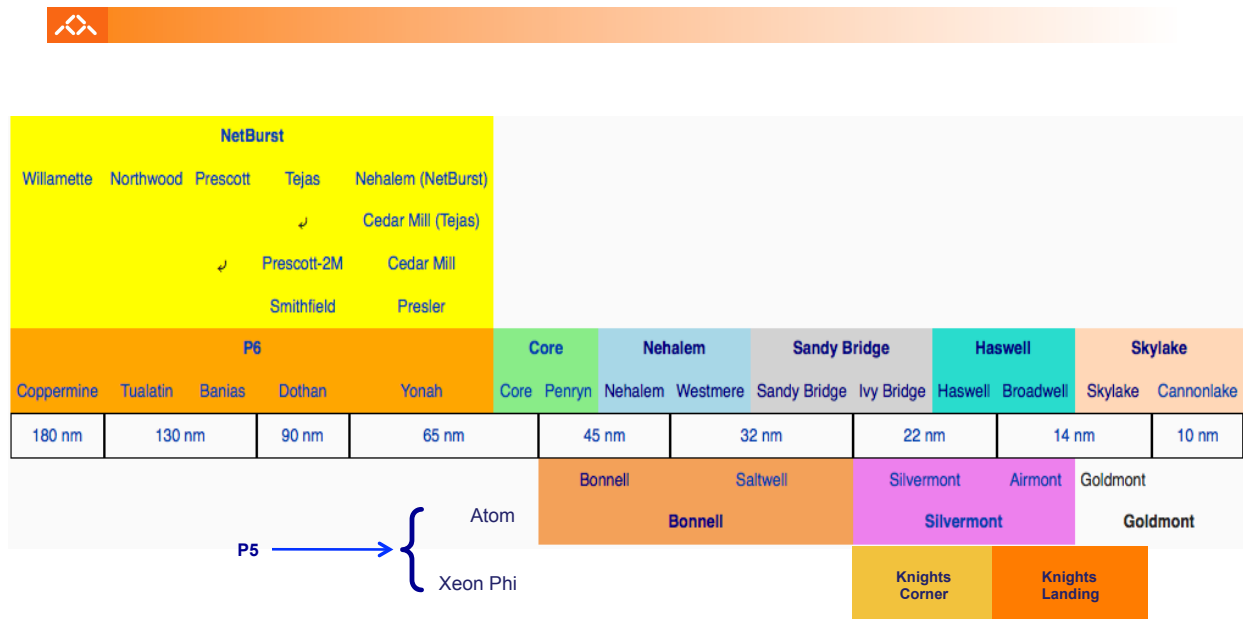**Three** Modes. Selected at boot

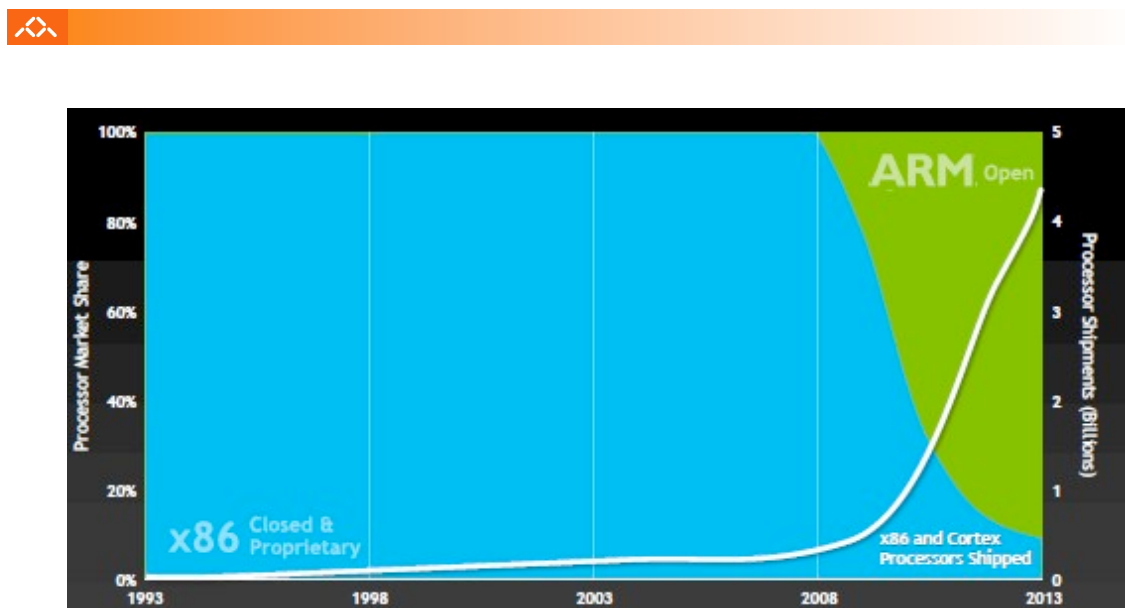Cache Mode      Flat Mode      Hybrid Mode

Omni-path not shown

Package

*AJProença, Parallel Computing, MiEI, UMinho, 2018*

# Internal x86 roadmap

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **NetBurst** | | | | | | | | | | |
| Willamette | Northwood | Prescott | Tejas | Nehalem (NetBurst) | | | | | | |
| | | | ↵ | Cedar Mill (Tejas) | | | | | | |
| | | ↵ | Prescott-2M | Cedar Mill | | | | | | |
| | | | Smithfield | Presler | | | | | | |
| **P6** | | | | | **Core** | **Nehalem** | | **Sandy Bridge** | | **Haswell** | | **Skylake** | |
| Coppermine | Tualatin | Banias | Dothan | Yonah | Core | Penryn | Nehalem | Westmere | Sandy Bridge | Ivy Bridge | Haswell | Broadwell | Skylake | Cannonlake |
| 180 nm | 130 nm | 90 nm | 65 nm | | 45 nm | | 32 nm | | 22 nm | | 14 nm | | 10 nm | |

| Atom | Bonnell | Saltwell | Silvermont | Airmont | Goldmont |
|---|---|---|---|---|---|
| | **Bonnell** | | **Silvermont** | | **Goldmont** |

**P5** → { Atom / Xeon Phi }
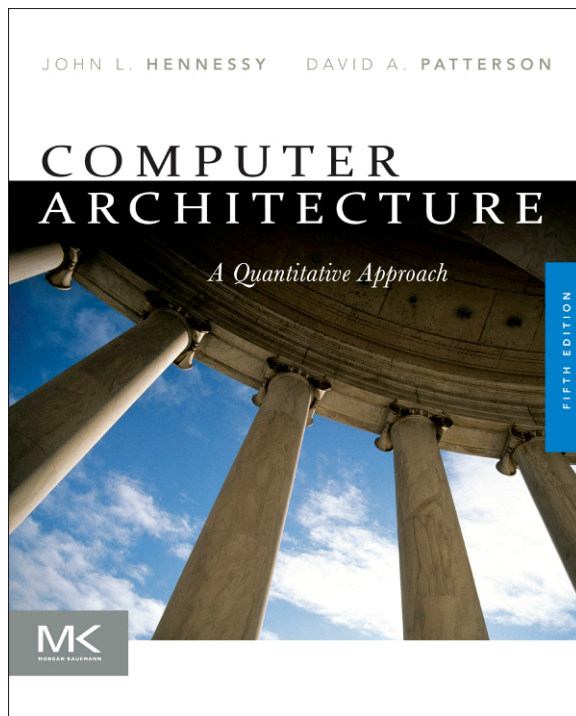
| Xeon Phi | Knights Corner | Knights Landing |
|---|---|---|

# Processadores Intel x86 versus ARM

# Key textbook in Advanced Architecture

**Computer Architecture, 5th Edition**

**Hennessy & Patterson**

## Table of Contents

# Recommended textbook (1)

## Table of Contents

### Programming Massively Parallel Processors

David B. Kirk
Wen-mei W. Hwu

SECOND EDITION

A Hands-on Approach

**Contents**

## Understanding Performance

- Algorithm + Data Structures
  - Determines number of operations executed
  - Determines how efficient data is assessed
- Programming language, compiler, architecture
  - Determine number of machine instructions executed per operation
- Processor and memory system
  - Determine how fast instructions are executed
- I/O system (including OS)
  - Determines how fast I/O operations are executed

## Response Time and Throughput

- Response time
  - How long it takes to do a task
- Throughput
  - Total work done per unit time
    - e.g., tasks/transactions/… per hour
- How are response time and throughput affected by
  - Replacing the processor with a faster version?
  - Adding more processors?
- We'll focus on response time for now…

## CPU Time (single-core)

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

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

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

$$\text{Clock Cycles} = \text{Instruction Count} \times \text{Cycles per Instruction}$$

$$\text{CPU Time} = \text{Instruction Count} \times \text{CPI} \times \text{Clock Cycle Time}$$

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

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

---

*Performance Summary* **(single-core)**

**The BIG Picture**

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

- Performance depends on
  - Algorithm: affects IC, possibly CPI
  - Programming language: affects IC, CPI
  - Compiler: affects IC, CPI
  - Instruction set architecture: affects IC, CPI, $T_c$
  - Processor design: ILP, memory hierarchy, ...

⬦

## The BIG Picture

- Pipelining improves performance by increasing instruction throughput
  - Executes multiple instructions in parallel
  - Each instruction has the same latency
- Subject to hazards
  - Structure, data, control
- Instruction set design affects complexity of pipeline implementation

---

⬦

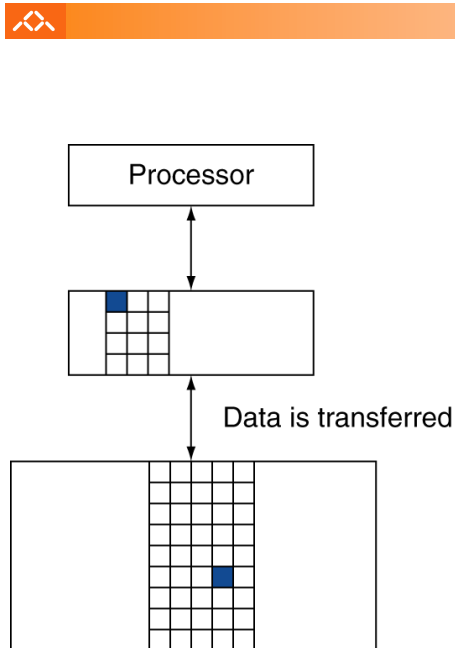## The BIG Picture

- 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

## *Memory Hierarchy Levels*

Processor

Data is transferred

- Block (aka line): unit of copying
  - May be multiple words
- If accessed data is present in upper level
  - Hit: access satisfied by upper level
    - Hit ratio: hits/accesses
- If accessed data is absent
  - Miss: block copied from lower level
    - Time taken: miss penalty
    - Miss ratio: misses/accesses = 1 – hit ratio
  - Then accessed data supplied from lower level

---

## *The Memory Hierarchy*

### The BIG Picture

- Common principles apply at all levels of the memory hierarchy
  - Based on notions of caching
- Decisions at each level in the hierarchy
  - Block placement
  - Finding a block
  - Replacement on a miss
  - Write policy

§5.5 A Common Framework for Memory Hierchies

- Primary cache private to CPU/core
  - Small, but fast
- Level-2 cache services misses from primary cache
  - Larger, slower, but still faster than main memory
- High-end systems include L3 cache
- Main memory services L2/3 cache misses

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

*AJProença, Parallel Computing, MiEI, UMinho, 2018/19*                                                                23

*Memory Hierarchy*

Introduction



(a) Memory hierarchy for server



(b) Memory hierarchy for a personal mobile device