



Estrutura do tema ISC

1. Representação de informação num computador
2. Organização e estrutura interna dum computador
3. Execução de programas num computador
4. Análise das instruções de um processador
5. Evolução da tecnologia e da eficiência



Componentes (físicos) a analisar:

- a unidade de processamento / o processador:
 - o nível ISA (*Instruction Set Architecture*): tipos/formatos de instruções, acesso a operandos, CISC/RISC...
 - **CISC versus RISC**
 - **melhoria de eficiência** no processador: com paralelismo
 - **melhoramentos** fora do processador (ou core)
 - **evolução** da arquitetura x86 da Intel
- a hierarquia de memória:
cache, memória virtual, ...
- periféricos...

CISC versus RISC

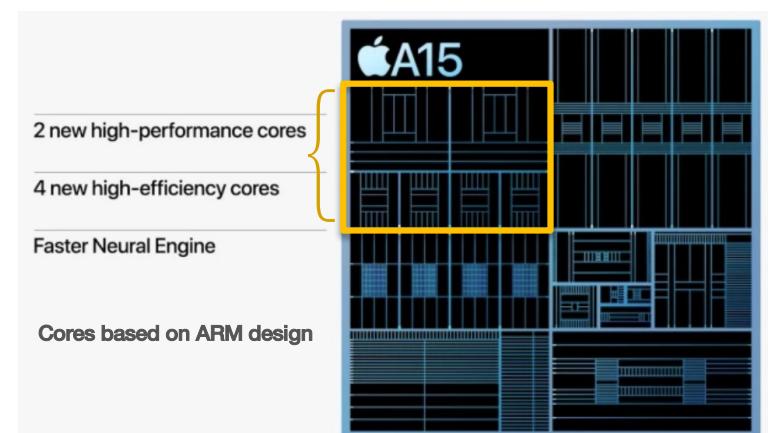


A “revolução” dos *Reduced Instruction Set Computers*

Caracterização das arquiteturas RISC

- conjunto reduzido e simples de instruções
- formatos simples de instruções
- uma operação elementar por ciclo máquina
- operandos sempre em registos
- modos simples de endereçamento à memória

Arquiteturas RISC: **em todos os smartphones!**



Eficiência nos Sistemas de Computação: oportunidades para melhorar



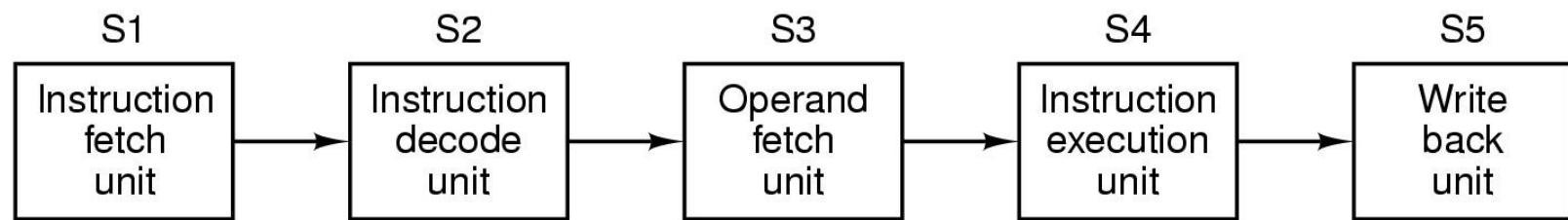
Oportunidades para melhorar o desempenho / eficiência

- com introdução de **parallelismo**
 - na execução de instruções em binário (*Instruction Level Parallelism*)
 - paralelismo desfasado ou execução encadeada (*pipeline*)
 - paralelismo nos dados (processamento vetorial, MMX/SSE/AVX...)
 - paralelismo nas operações (VLIW, superescalar)
 - com execução de instruções fora de ordem (*out-of-order execution*)
 - no acesso à memória
 - paralelismo desfasado (*interleaving*)
 - paralelismo "real" (maior largura do bus)
 - ao nível da aplicação (sistemas concorrentes/paralelos/distribuídos)
 - com fios de execução (*multithreading* => *multicore/multichip* c/ mem partilhada)
 - com processos (com memória distribuída)
- com **hierarquia de memória** (para diminuir latência)
 - *cache* ...

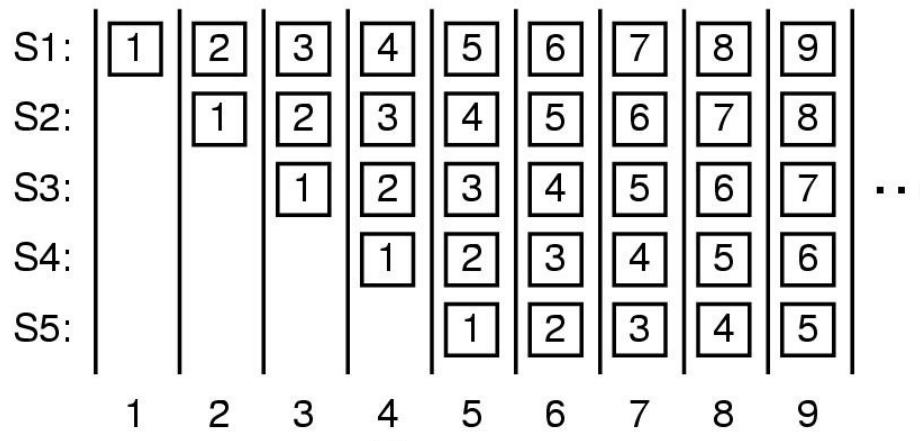
Paralelismo no processador: exemplo



Exemplo de *pipeline*

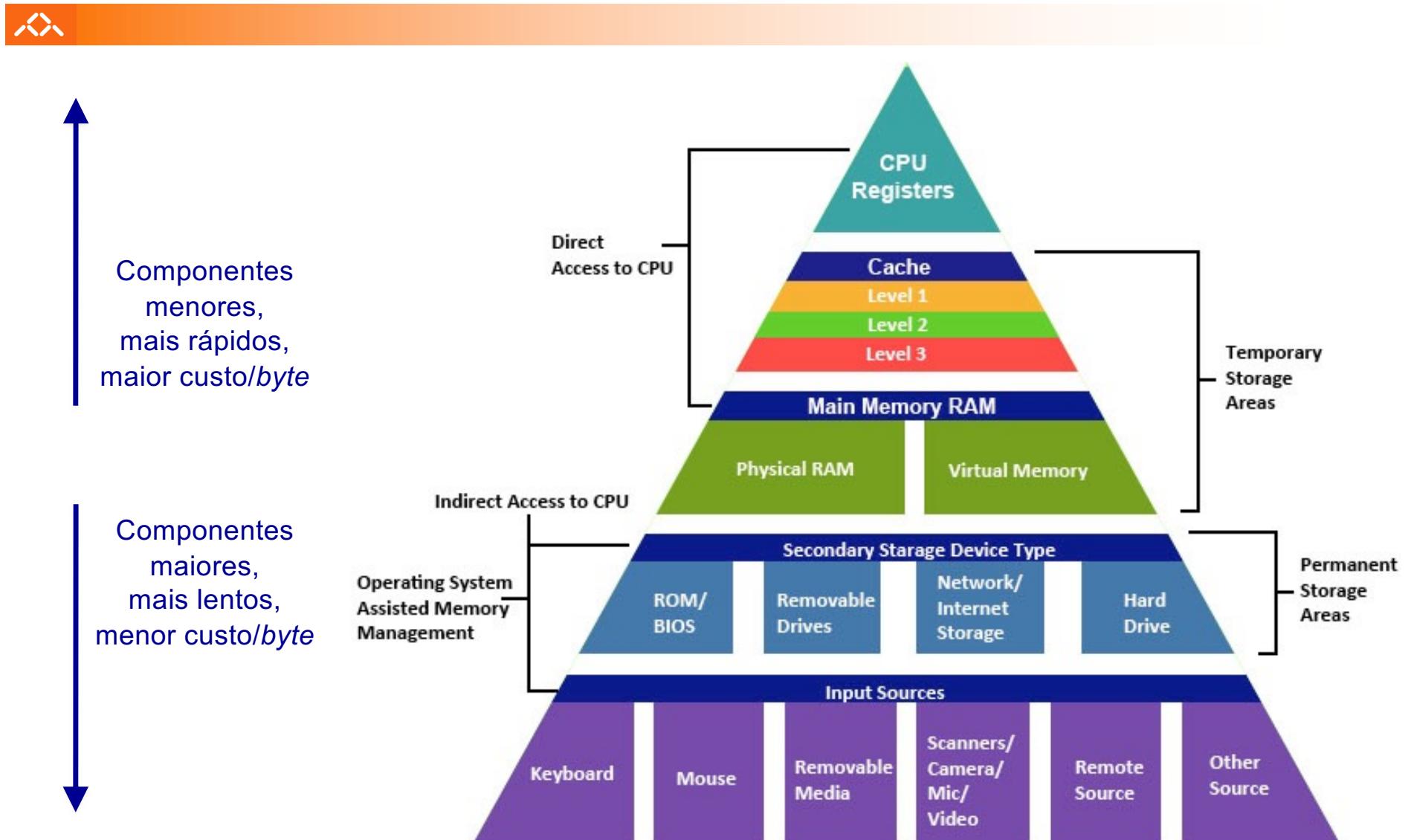


(a)



(b)

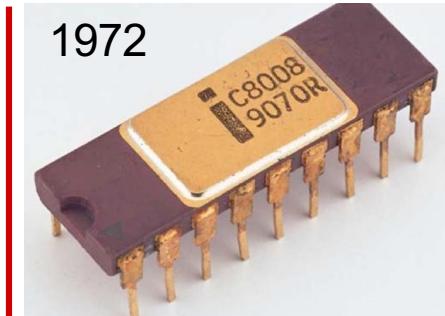
Hierarquia de memória



Evolução dos processadores da Intel até à família Intel x86



4004: 1º processador num único *chip* (microprocessador)



8008 e 8080: 1ºs microprocessadores de 8 bits



8088 e 8086: 1ºs microprocessadores de 16 bits (selecionados para o IBM PC)



Evolução do Intel x86 : pré-Pentium (visão do programador) (1)



<i>Nome</i>	<i>Data</i>	<i>Nº transístores</i>	
8086	1978	29K	<ul style="list-style-type: none">– processador 16-bits (<i>registos + ALU</i>); base do <i>IBM PC & DOS</i>– espaço de endereçamento limitado a <i>1MiB</i> (<i>DOS apenas vê 640Ki</i>)
80286	1982	134K	<ul style="list-style-type: none">– endereço 24-bits e <i>protected-mode</i>; base do <i>IBM PC-AT & Windows</i>
386	1985	275K	→ <u>primeiro IA-32 !!</u>
			<ul style="list-style-type: none">– estendido para 32-bits: <i>registos + op. inteiros + endereçamento</i>– memória segmentada+paginada, capaz de correr <i>Unix</i>
486	1989	1.9M	<ul style="list-style-type: none">– integração num único chip: 386, co-proc 387, até 16KiB cache L1– poucas alterações na arquitetura interna do processador

Evolução do IA-32: família Pentium *(visão do programador) (2)*

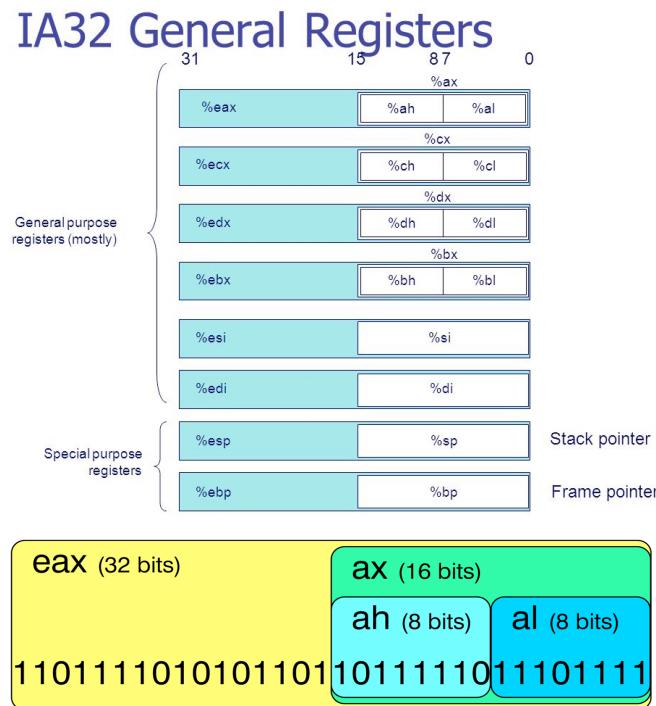


Pentium	1993	3.1M	(= P5 , aka i586)
			– arquitetura superescalar, com 2 pipelines de inteiros (de 5 níveis)
PentiumPro	1995	5.5M	(= P6 , aka i686)
			– <i>out-of-order execution, 14 níveis pipeline, 3-issue superscalar</i>
			– <i>endereço 36-bits, cache L2 on-package</i>
Pentium/MMX	1997	4.5M	
			– SIMD: opera com vetores de 64-bits, tipo <i>int</i> de 1, 2, ou 4 bytes
Pentium II	1997	7.5M	(= Pro + MMX)
Pentium III	1999	8.2M	(+Celeron, +Xeon)
			– “Streaming SIMD Ext”, SSE: vetores 128-bits, <i>int/fp</i> 1/2/4 bytes
Pentium 4	2000	42M	(= NetBurst , aka i786)
			– <i>trace cache, pipeline muito longo (20 ou 31 níveis), suporta multi-threading</i>
			– SSE2: <i>mais instruções e com dados fp de 8-bytes</i>
Pentium M	2003	77M	(= P-M)
			– arquitetura mais próxima do Pentium III (eficiência energética)

Evolução do IA-32 para Intel 64 (visão do programador) (3)



- IA-32 ou x86 open architecture cresce para 64-bits
 - HP e Intel propõem arquitetura incompatível c/ IA-32: IA-64 (Itanium PU)
 - AMD anuncia em 1999 extensão do x86: x86-64
 - Intel segue AMD: IA-32e (Fev-04), EM64T (Mar-04), ou Intel 64 (2006)
 - AMD64 e Intel 64 diferentes, mas compiladores usam sub-set comum



x86-64 Integer Registers

%rax	%eax
%rbx	%ebx
%rcx	%ecx
%rdx	%edx
%rsi	%esi
%rdi	%edi
%rsp	%esp
%rbp	%ebp
%r8	%r8d
%r9	%r9d
%r10	%r10d
%r11	%r11d
%r12	%r12d
%r13	%r13d
%r14	%r14d
%r15	%r15d

- Twice the number of registers
- Accessible as 8, 16, 32, 64 bits

Intel 64 ≠ IA-64 (Itanium): nos registros

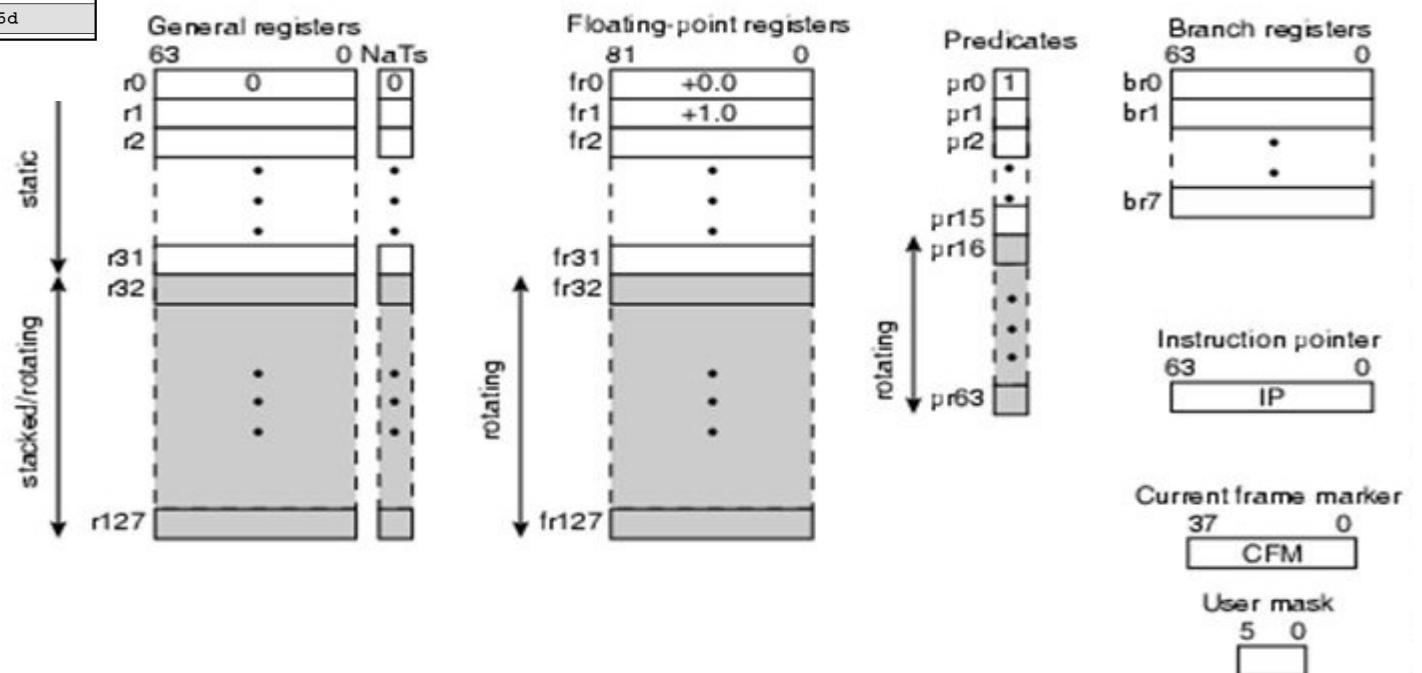


x86-64 Integer Registers

%rax	%eax
%rbx	%ebx
%rcx	%ecx
%rdx	%edx
%rsi	%esi
%rdi	%edi
%rsp	%esp
%rbp	%ebp
%r8	%r8d
%r9	%r9d
%r10	%r10d
%r11	%r11d
%r12	%r12d
%r13	%r13d
%r14	%r14d
%r15	%r15d

- Twice the number of registers
- Accessible as 8, 16, 32, 64 bits

IA-64 Register Set



Arquiteturas Intel 64 com maior integração (visão do programador) (1)

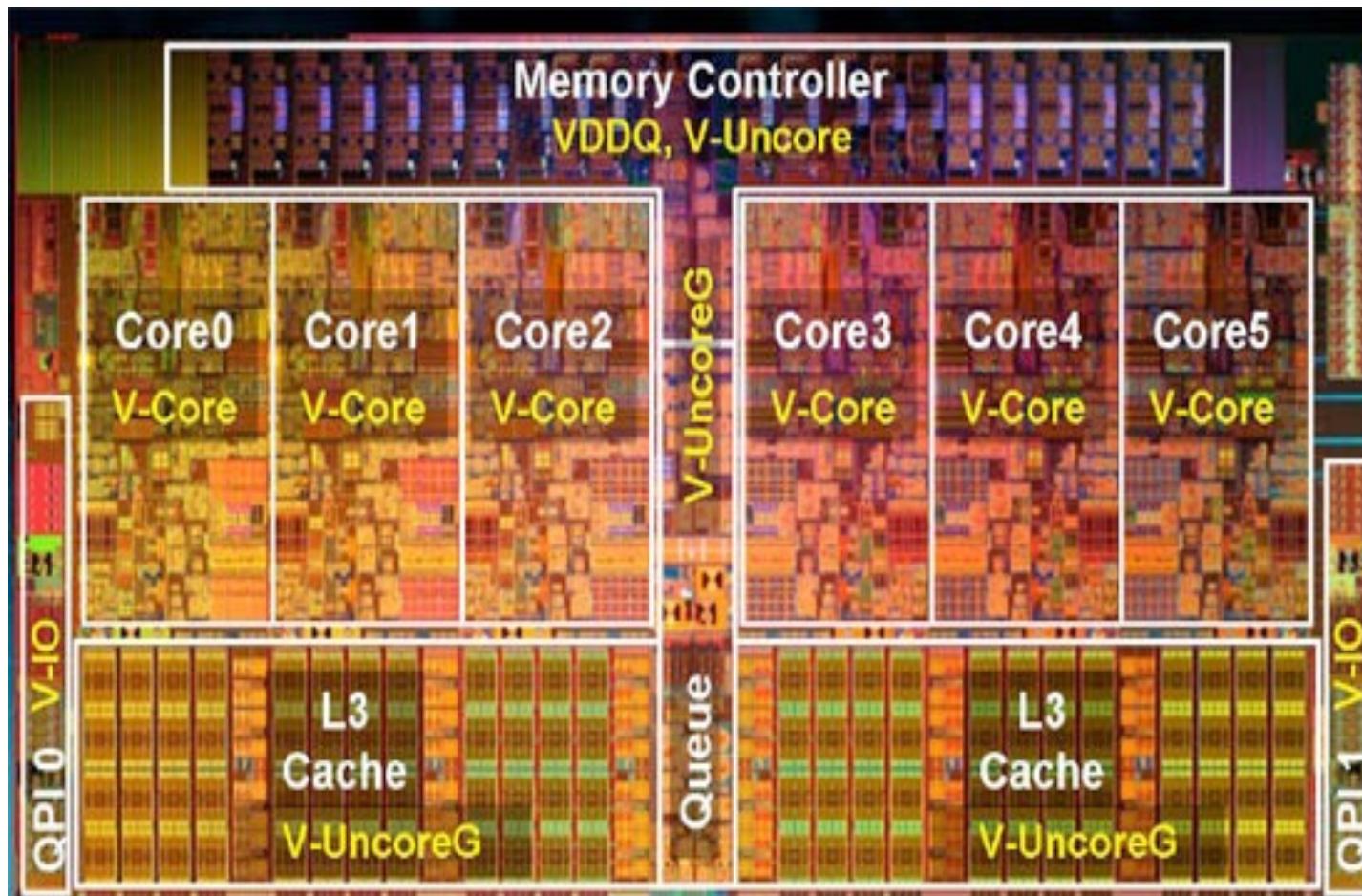


- arquitetura **Core** surge em 2006 (*151M transístores*)
 - desenvolvida pela mesma equipa que o P-M (Israel)
 - 14 níveis de pipeline (como P6), mas 4-issue superscalar
 - 2 níveis de cache on-chip
 - multi-core on-chip e virtualização por h/w
 - suporta fusão de instruções RISC (μ -ops na terminologia Intel)
 - arquitetura Core 2 é integralmente 64-bit (Intel 64)
- arquitetura **Nehalem** anunciada em 2008 (*731M transístores*)
 - inspirada no NetBurst (com multi-threading e maiores frequências de clock)
 - 2 a 8 cores por chip, com cache L3 on-chip
 - com conexão ponto-a-ponto inter-CPU_chips
 - integra controlador de memória, rumo a servidores com arquitetura NUMA

Xeon Nehalem



Intel Hex-Core Nehalem (1.17 mil milhões de transístores)



Arquiteturas Intel 64 com maior integração

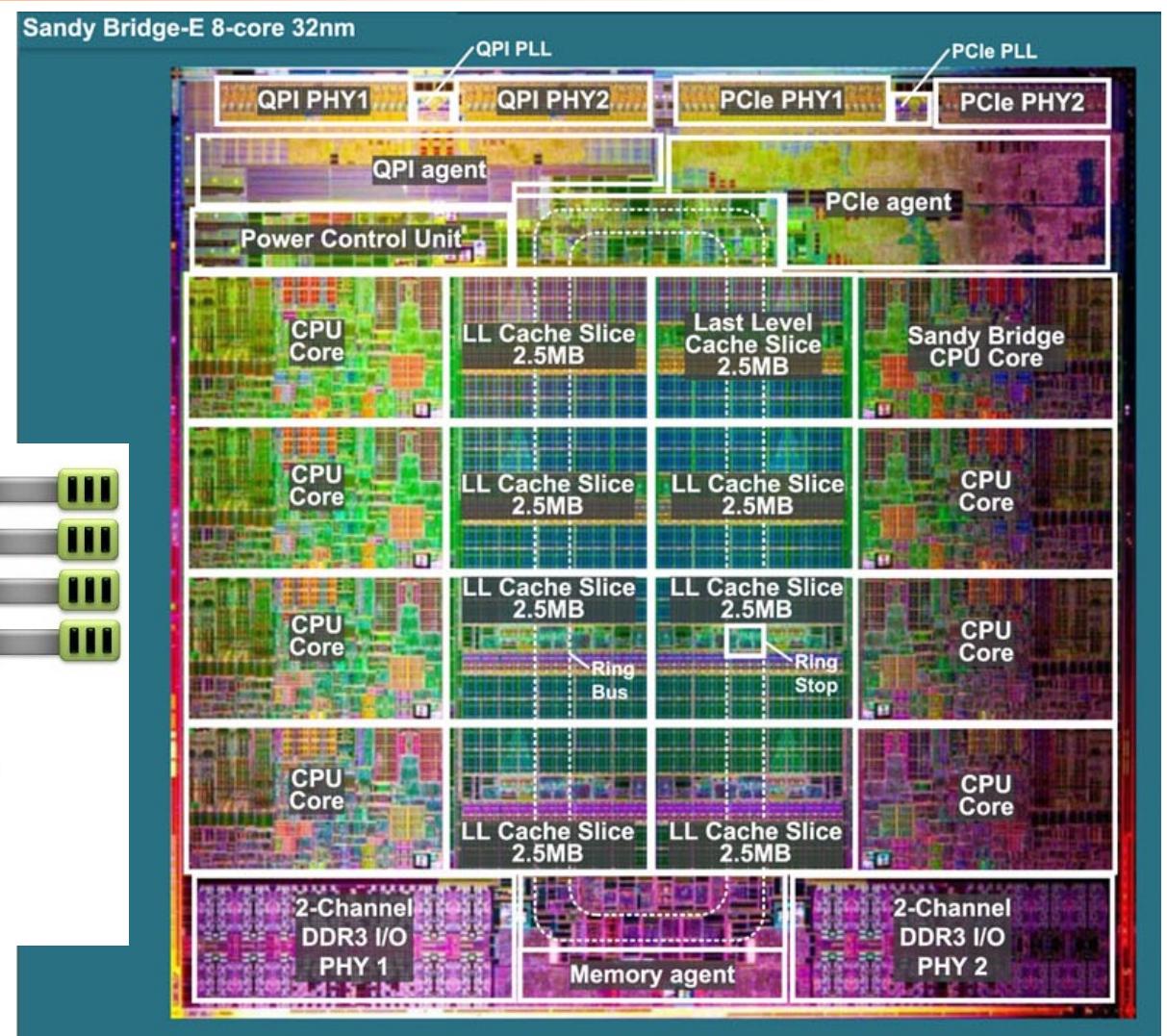
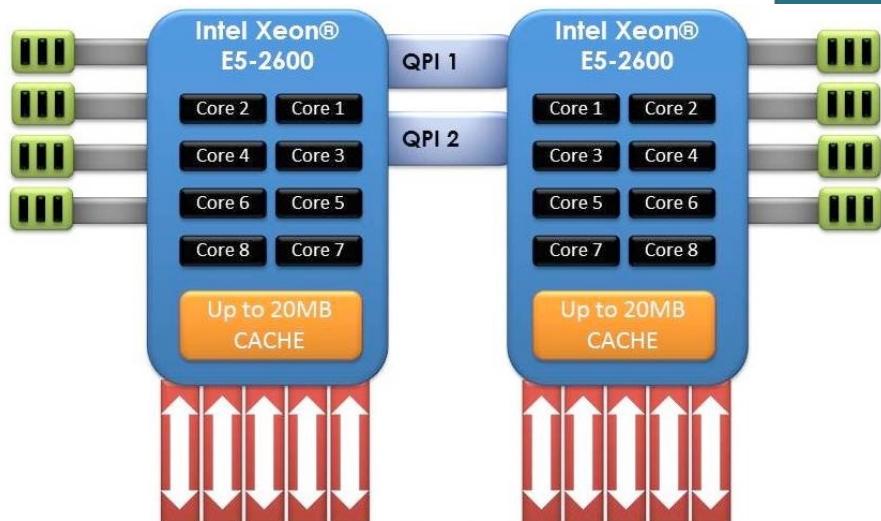
(visão do programador) (2)



- arquitetura **Sandy Bridge** anunciada em 2010 ($>1,000M$ transístores)
 - *ensaio de arquitetura híbrida multicore, integrando o processador gráfico*
 - *interface com o processador gráfico através da cache L3*
 - *processamento vetorial de fp estendido para 256-bits (AVX)*
 - *integração no chip do interface PCIe 16x*
- arquitetura **Haswell** anunciada em 2013 ($>5,500M$ transístores)
 - *nível adicional de cache para μ -ops (formato RISC)*
 - *processamento vetorial integral com 256-bits (AVX2)*
 - *2 unidades vetoriais para operações com inteiros*
 - *até 22 cores e 55 MiB de cache L3 (Xeon)*
 - *introdução de cache L4 de 128 MiB , eDRAM (off-chip. on-package)*
- arquitetura **Skylake** anunciada em 2015 ($>8,000M$ transístores)
 - *mais 1 unidade vetorial para operações com inteiros (total: 3)*
 - *AVX-512 para topo de gama Xeon (Skylake-SP)*

Xeon Sandy Bridge

Arquitetura NUMA

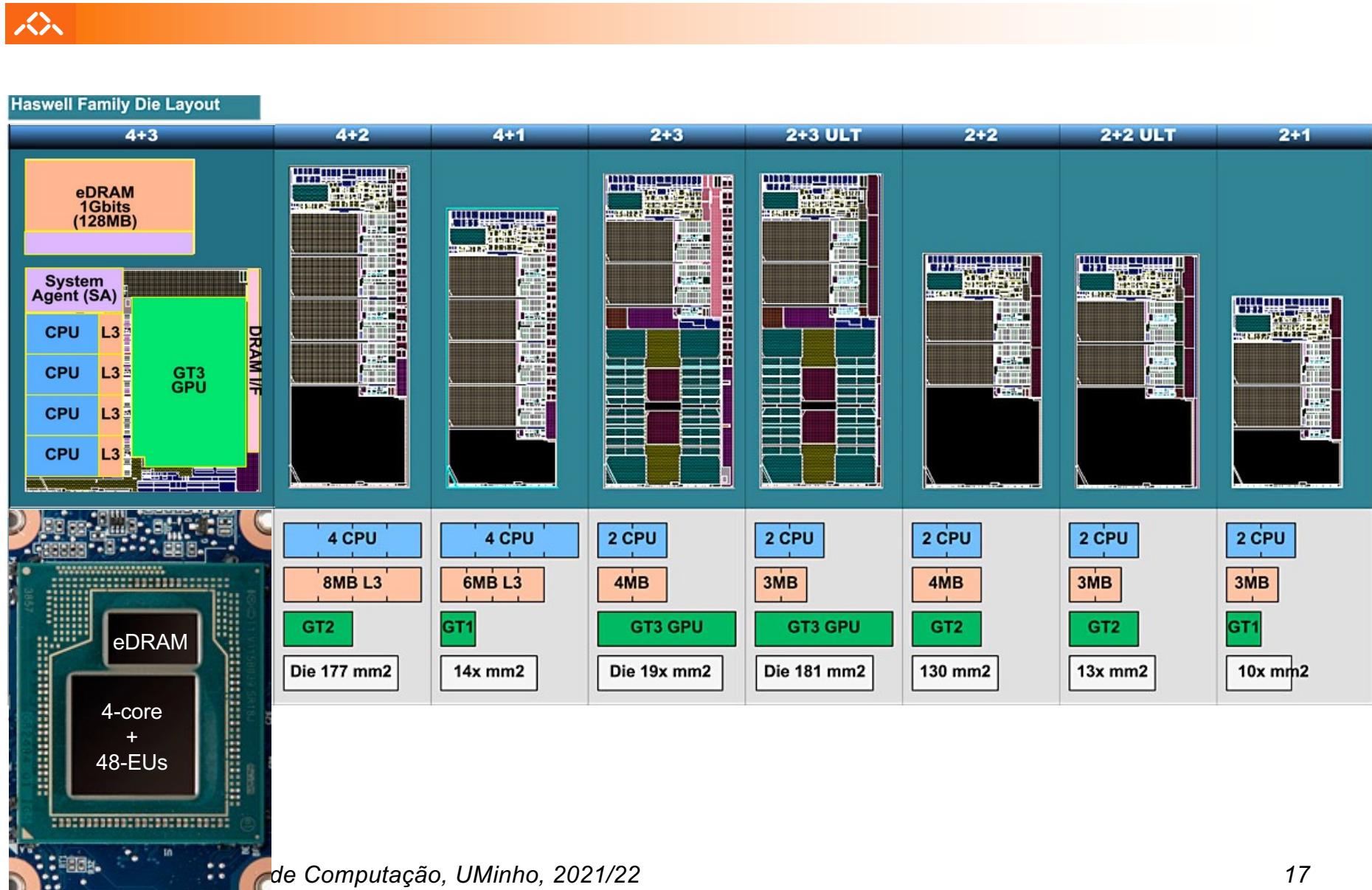


Arquiteturas Intel 64 com maior integração (visão do programador) (3)

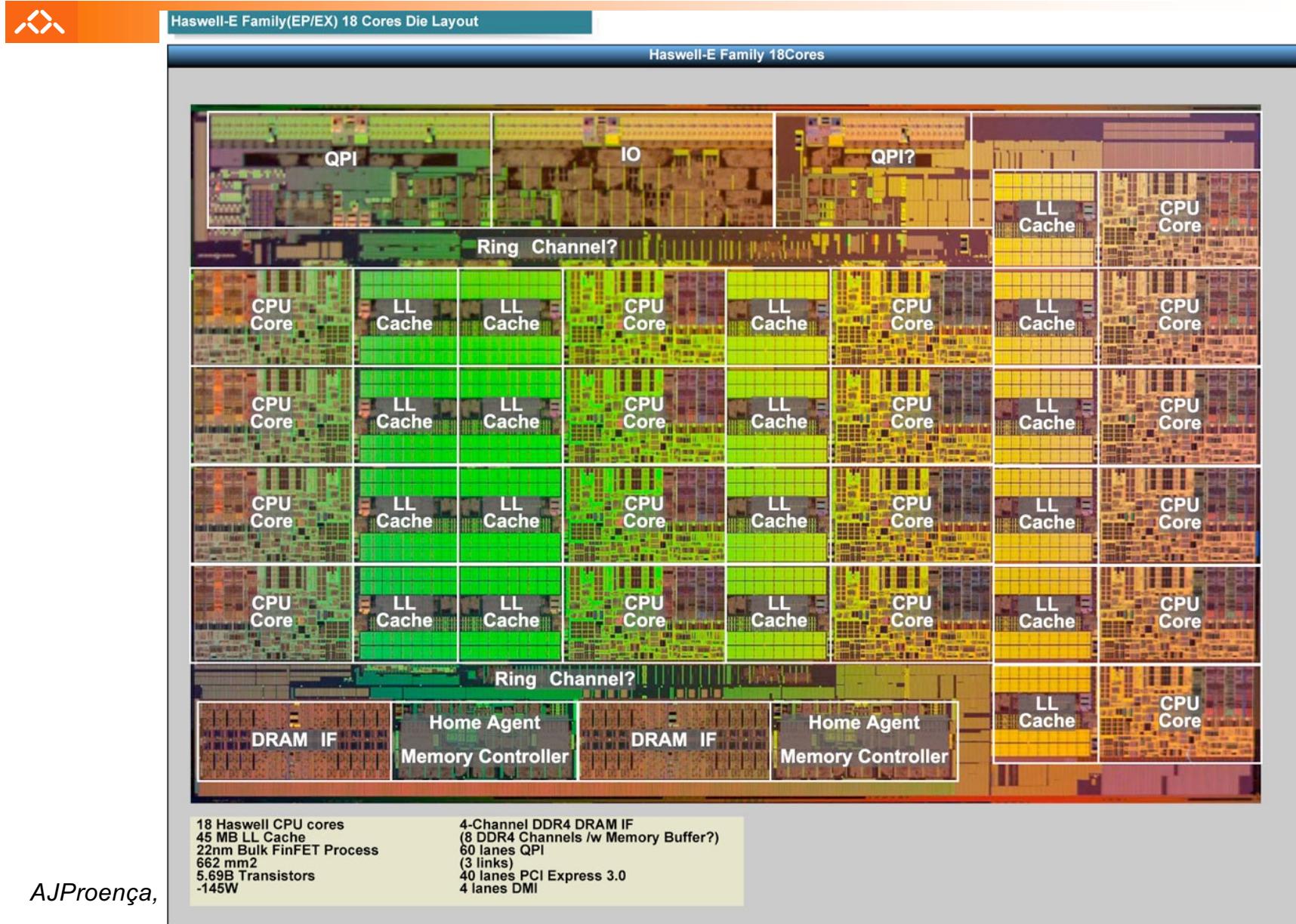


- arquitetura **Sandy Bridge** anunciada em 2010 ($>1,000M$ transístores)
 - *ensaio de arquitetura híbrida multicore, integrando o processador gráfico*
 - *interface com o processador gráfico através da cache L3*
 - *processamento vetorial de fp estendido para 256-bits (AVX)*
 - *integração no chip do interface PCIe 16x*
- arquitetura **Haswell** anunciada em 2013 ($>5,500M$ transístores)
 - *nível adicional de cache para μ -ops (formato RISC)*
 - *processamento vetorial integral com 256-bits (AVX2)*
 - *2 unidades vetoriais para operações com inteiros*
 - *até 22 cores e 55 MiB de cache L3 (Xeon)*
 - *introdução de cache L4 de 128 MiB , eDRAM (off-chip, on-package)*
- arquitetura **Skylake** anunciada em 2015 ($>8,000M$ transístores)
 - *mais 1 unidade vetorial para operações com inteiros (total: 3)*
 - *AVX-512 para topo de gama Xeon (Skylake-SP)*

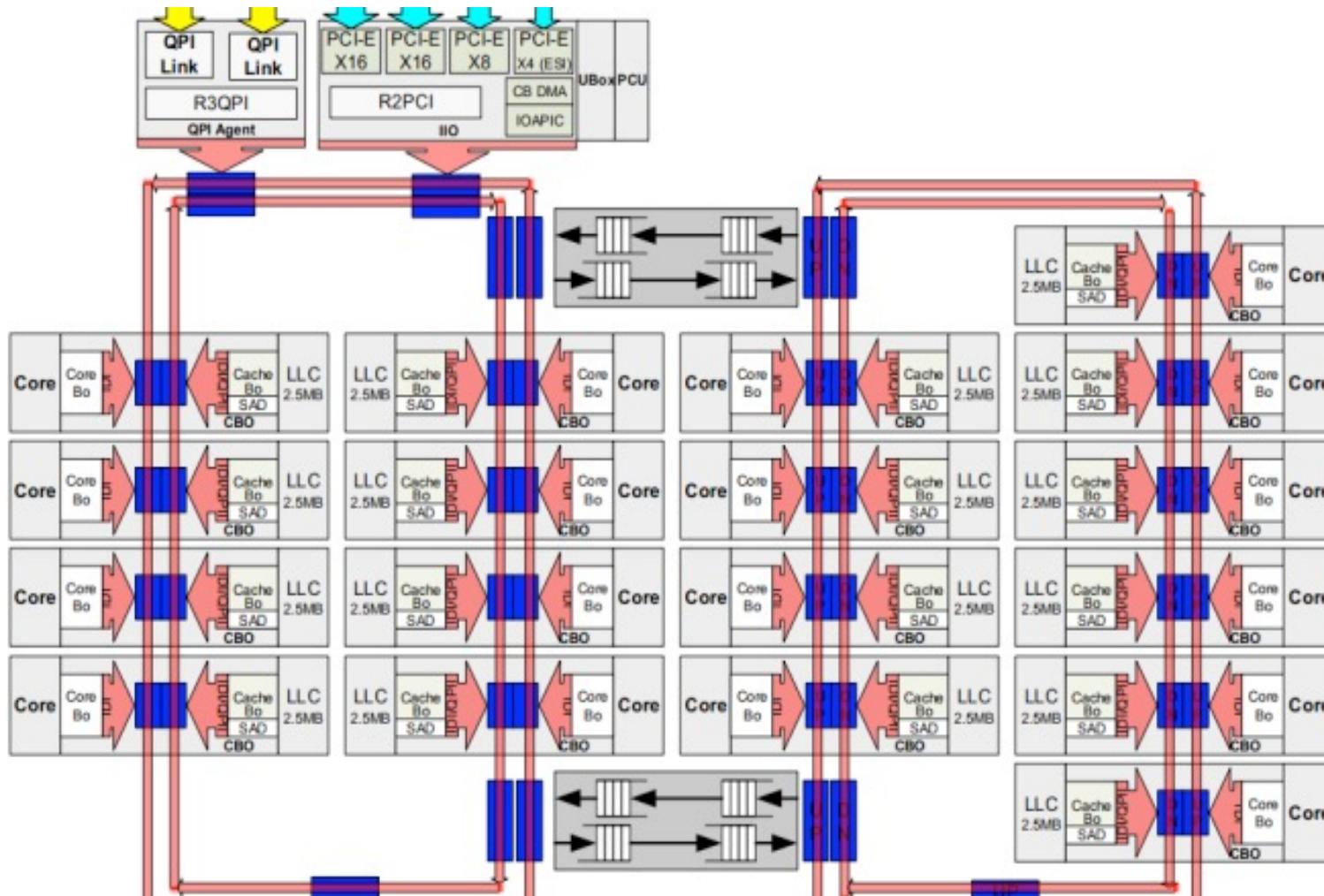
Lançamento da Intel em 2013: 8 configurações para i5 & i7 Haswell



Lançamento da Intel em 2016: 18-core Xeon Haswell



Intel 18-core Xeon Haswell



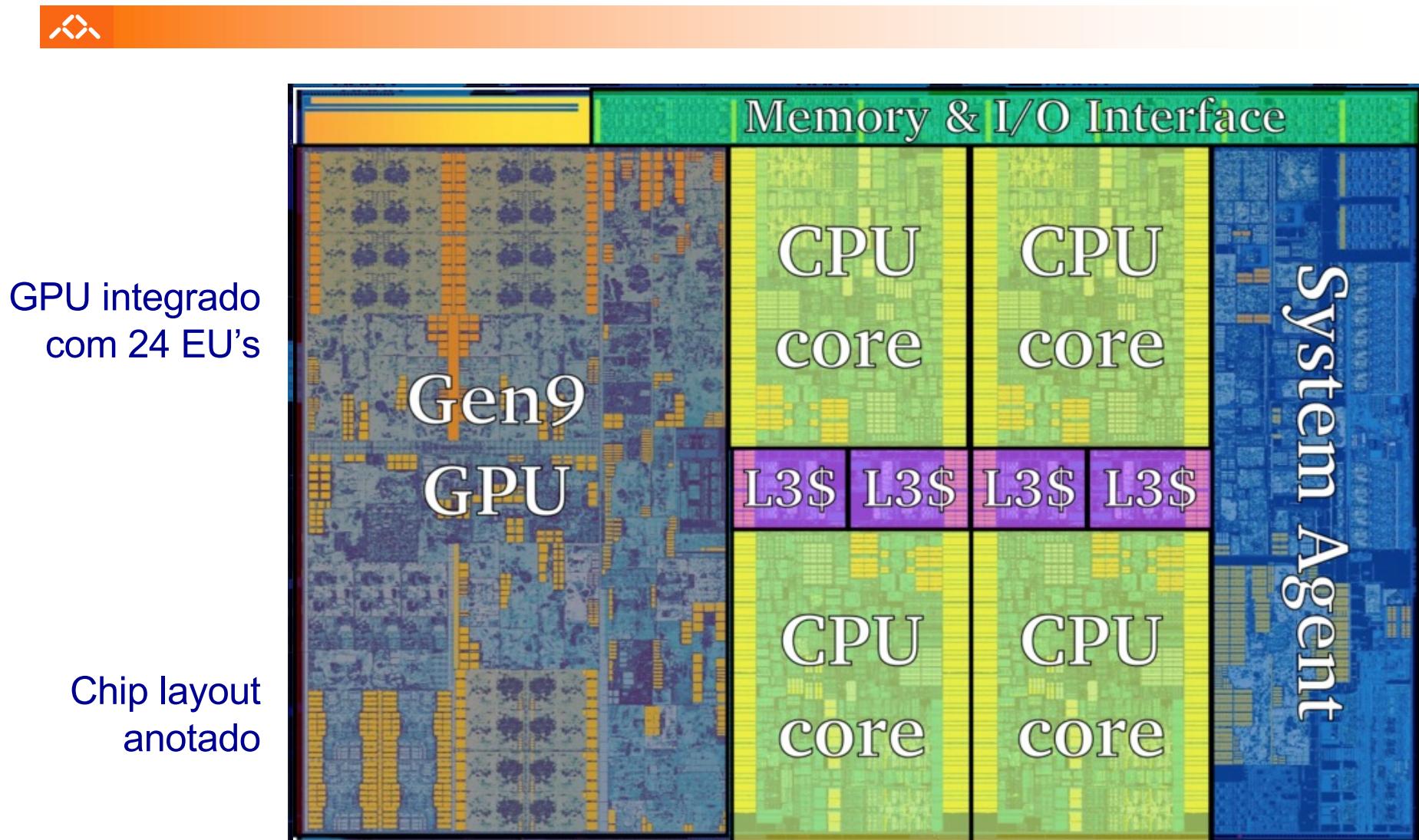
Arquiteturas Intel 64 com maior integração

(visão do programador) (4)

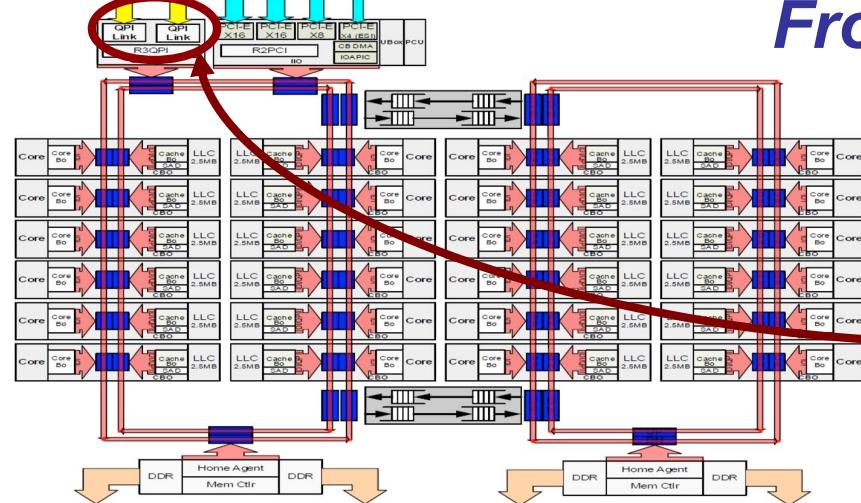


- arquitetura **Sandy Bridge** anunciada em 2010 ($>1,000M$ transístores)
 - ensaio de arquitetura híbrida multicore, integrando o processador gráfico
 - interface com o processador gráfico através da cache L3
 - processamento vetorial de fp estendido para 256-bits (AVX)
 - integração no chip do interface PCIe 16x
- arquitetura **Haswell** anunciada em 2013 ($>5,500M$ transístores)
 - nível adicional de cache para μ -ops (formato RISC)
 - processamento vetorial integral com 256-bits (AVX2)
 - 2 unidades vetoriais para operações com inteiros
 - até 22 cores e 55 MiB de cache L3 (Xeon)
 - introdução de cache L4 de 128 MiB , eDRAM (off-chip, on-package)
- arquitetura **Skylake** anunciada em 2015 ($>8,000M$ transístores)
 - mais 1 unidade vetorial para operações com inteiros (total: 3)
 - AVX-512 para topo de gama Xeon (Skylake-SP)
- arquitetura Xeon a seguir: **Sunny Cove**, anunciada a 6-abr-21
 - Xeon Ice Lake, redução de 14nm para 10nm, com 2 anos de atraso...
 - caches de maior dimensão

Intel 4-core Skylake



Chip layout
anotado



Broadwell

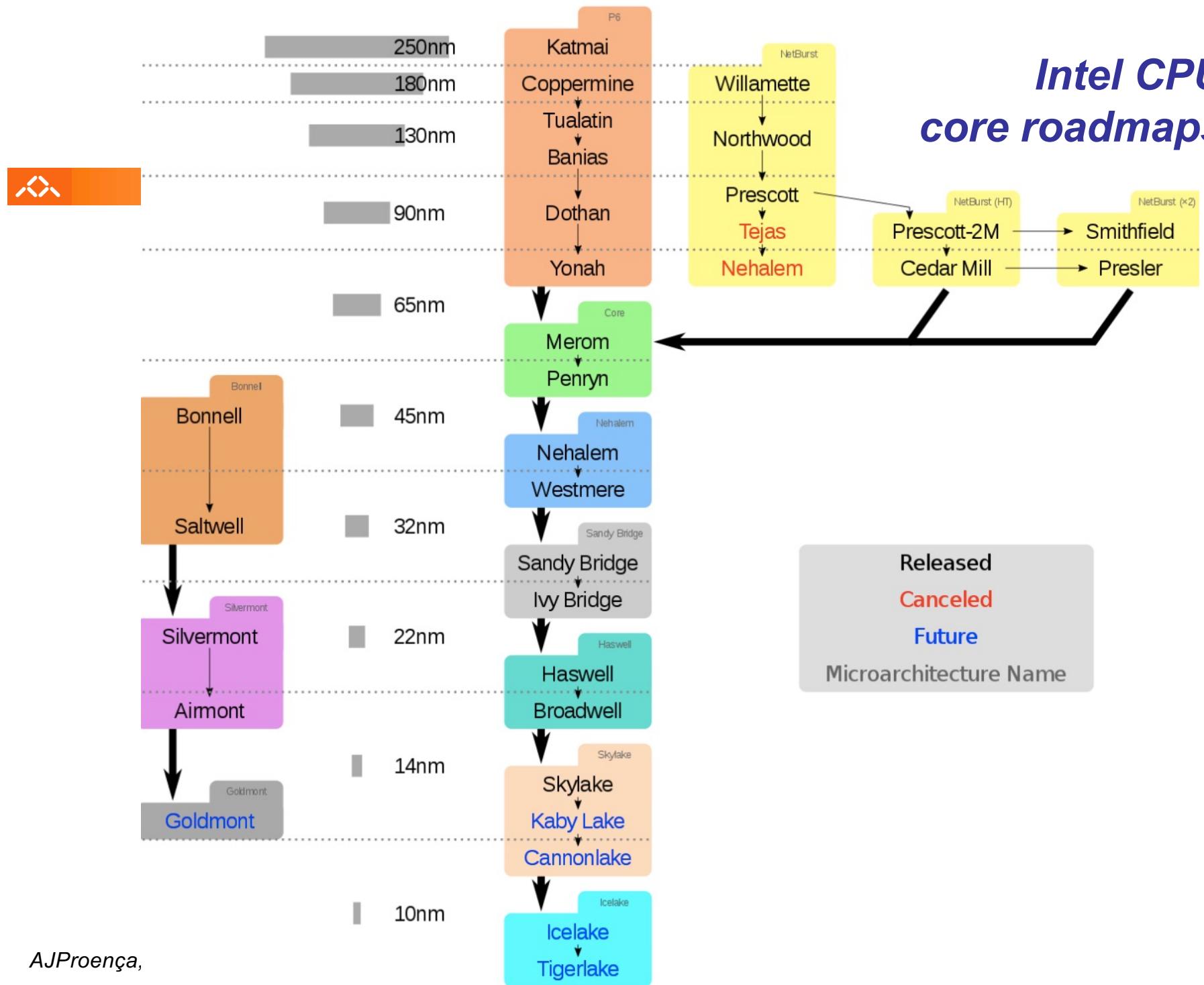
ring interconnection
does not scale for
large #cores

From Broadwell to Skylake (server): the move from ring to mesh

UPI required for dual-socket
(Ultra Path Interconnect)



Intel CPU core roadmaps



Novo modelo da Intel no desenvolvimento de processadores



- Segmentos de mercado: *desktop, mobile, server*
- Modelo Tick-Tock foi substituído em 2016 por novo modelo:

Process-Architecture-Optimization model

From Wikipedia, the free encyclopedia

Process–architecture–optimization is a processor development model adopted in 2016 by [Intel](#). Under this three-phase model, every [die shrink](#) is followed by a [microarchitecture](#) change and then by an optimization. It replaced the two-phase [Tick–tock model](#), adopted by Intel in 2006, because according to Intel the previous model was [and still is] no longer sustainable.^{[1][2][3][4]}

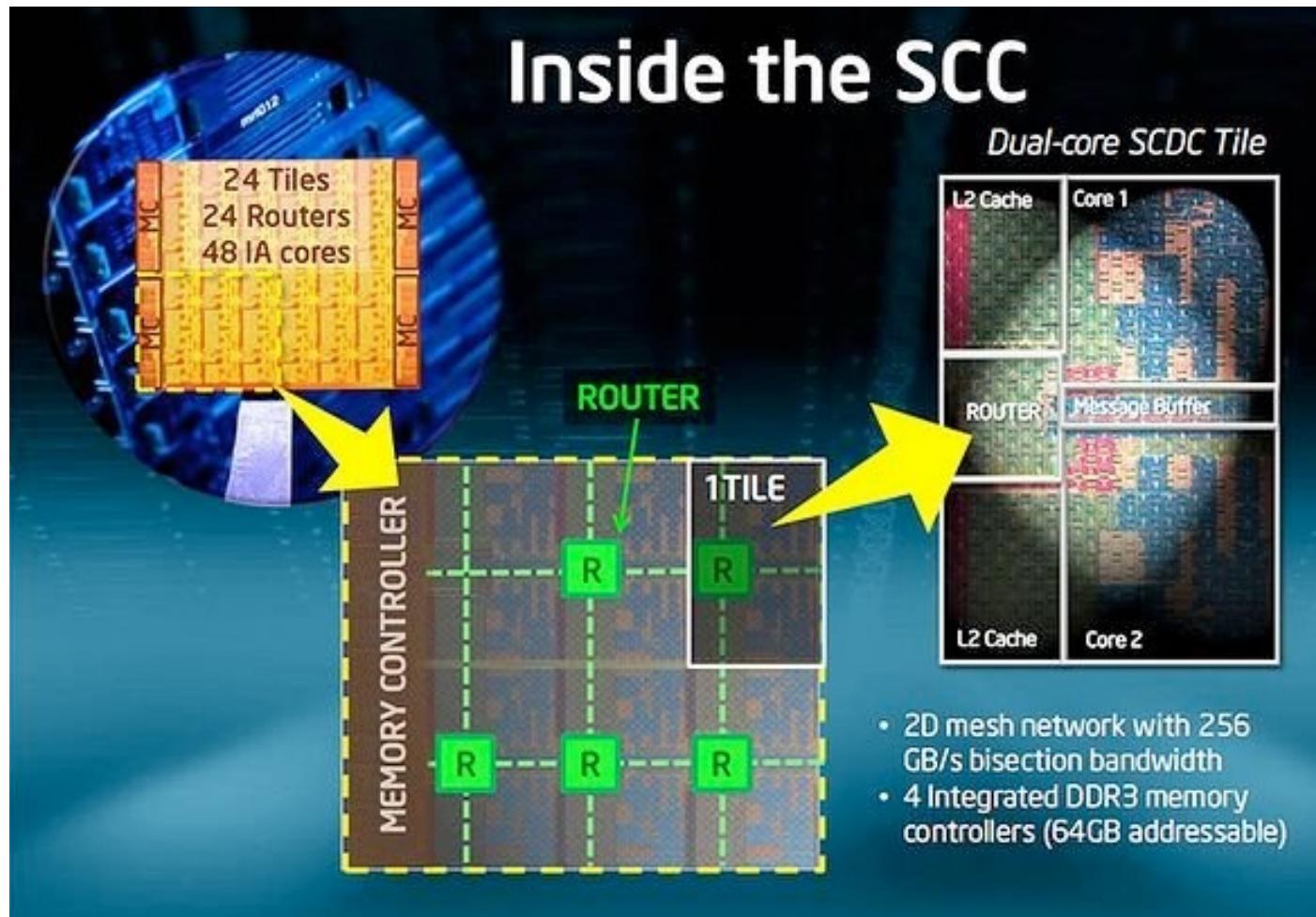
Scaling	Process	Architecture	Optimizations
14 nm	Broadwell	Skylake	Mainstream Desktop: Kaby Lake , Coffee Lake , Coffee Lake Refresh , Comet Lake , Rocket Lake Mobile: Kaby Lake , Kaby Lake Refresh , Coffee Lake , Whiskey Lake , Amber Lake , Coffee Lake Refresh , Comet Lake , Rocket Lake Workstation/Server: Kaby Lake , Cascade Lake , Cooper Lake
10 nm	Cannon Lake (mobile only)	Ice Lake (mobile + server)	Tiger Lake

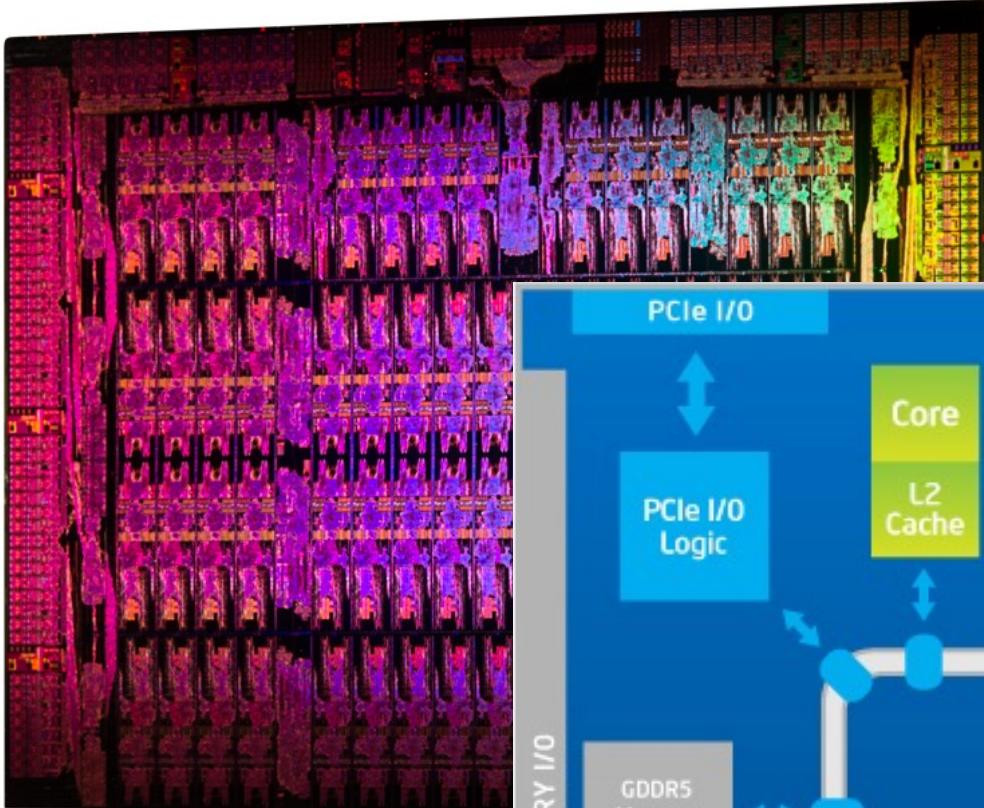
Roadmap dos Xeons de últimas gerações



Fabrication process	Micro-architecture	Code names	Core i generation	Release date	Processors		
					Enthusiast/WS	2P Server/WS	4P/8P Server
	Skylake	Skylake	6	2015-08-05 [33]	Skylake-X ^[34] Skylake-W	Skylake-SP (formerly Skylake-EP/-EX) ^[35]	
14 nm	Skylake + DLBoost	Cascade Lake	N/A	2019-04-02	Cascade Lake-X Cascade Lake-W Cascade Lake-SP		Cascade Lake-SP
	Skylake + DLBoost	Cooper Lake	N/A	2020	9]		Cooper Lake-SP
10 nm	Sunny Cove ^[41]	Ice Lake	10	2019-09 ^[a]	N/A	Ice Lake-SP ^[43]	
	Willow Cove ^[47]	Sapphire Rapids	N/A	2021 ^[48]		Sapphire Rapids-SP	

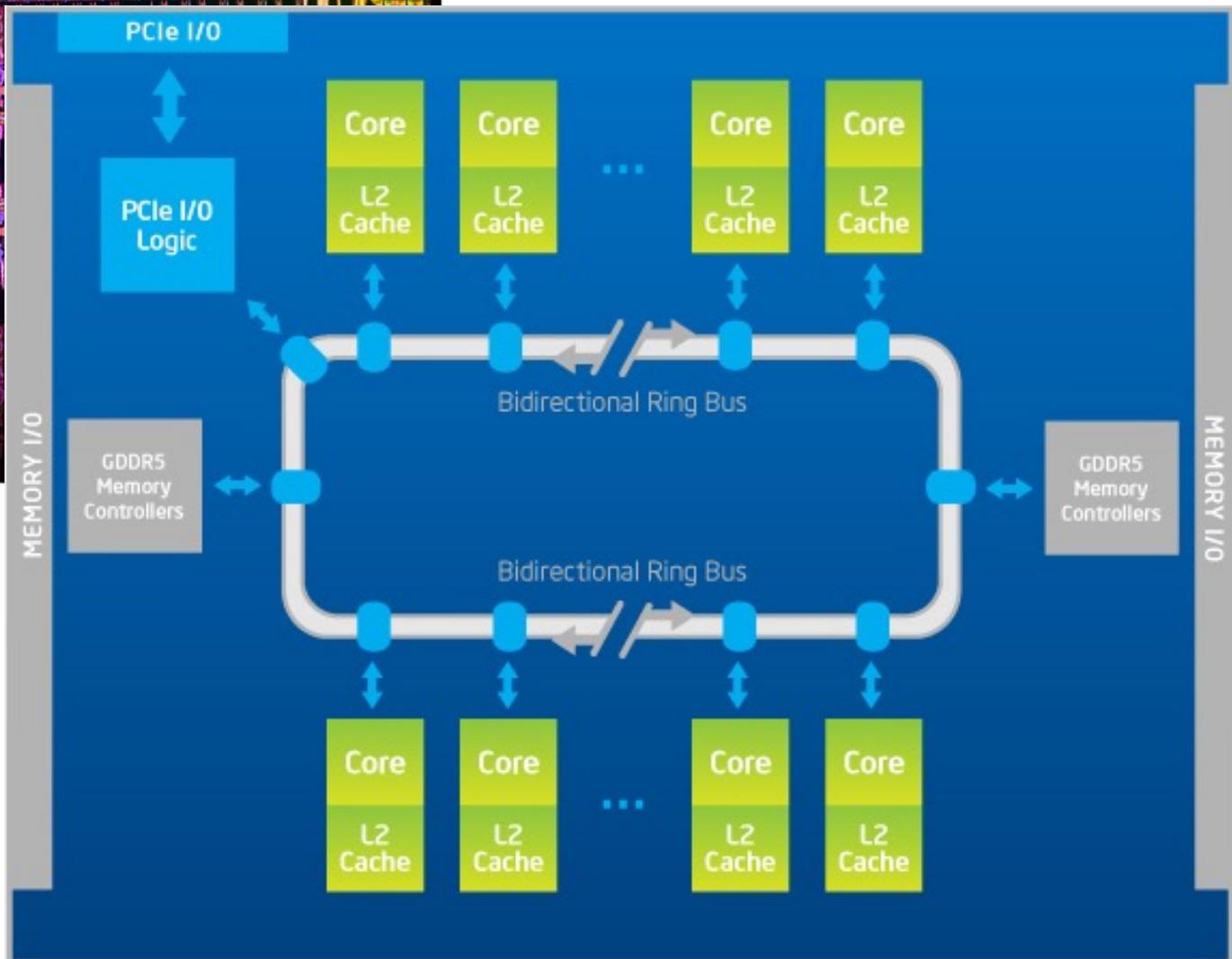
Protótipos da Intel em 2010/11: Single-chip Cloud Computer





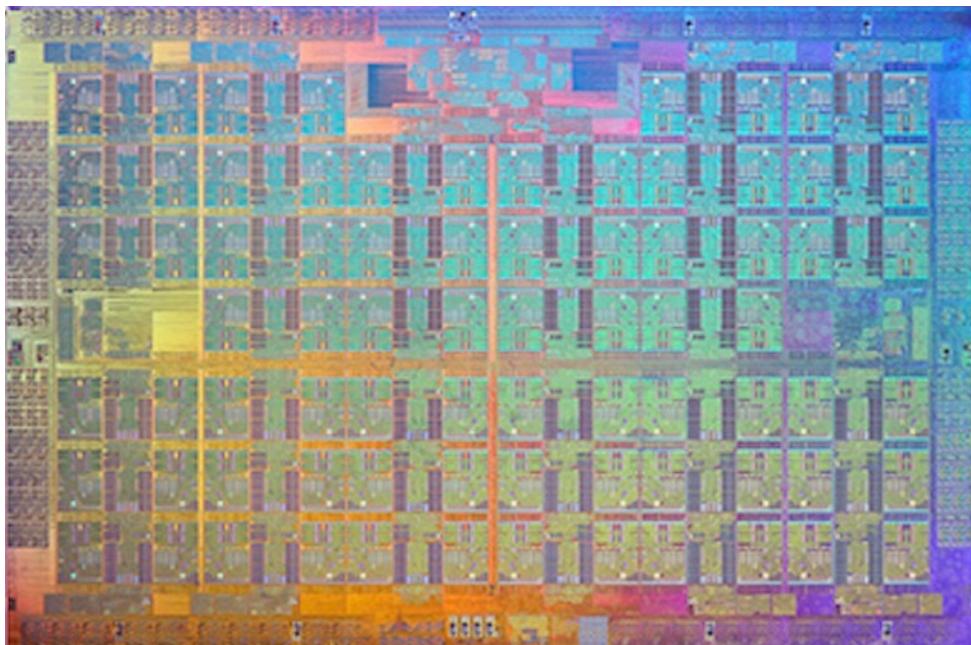
Xeon Phi
Knights Corner
KNC
(co-processador)

*Chips da Intel em 2012/13:
Xeon Phi com 60 cores
(apenas como co-processador)*

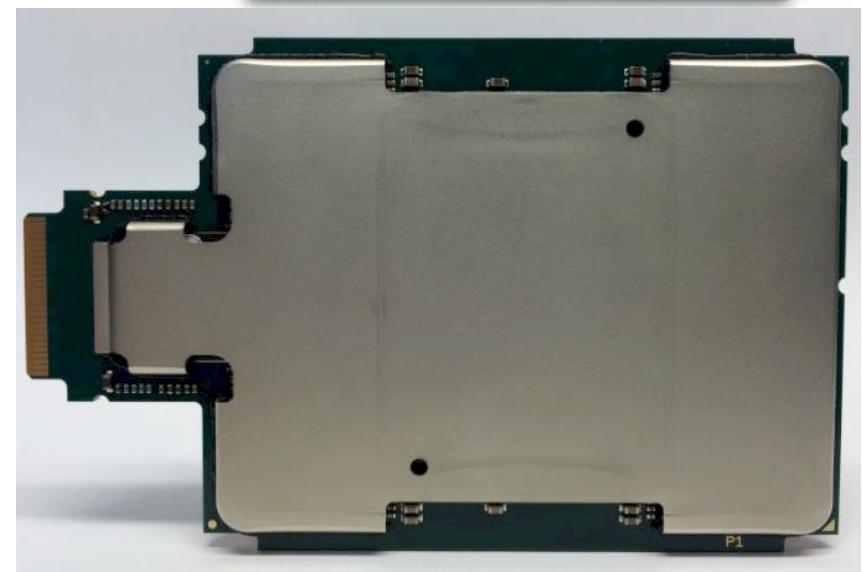
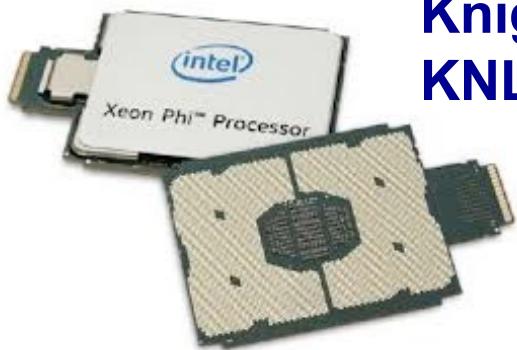


Chips da Intel em 2016: Xeon Phi até 72 cores

(como processador ou co-processador)



**Knights Landing
KNL**



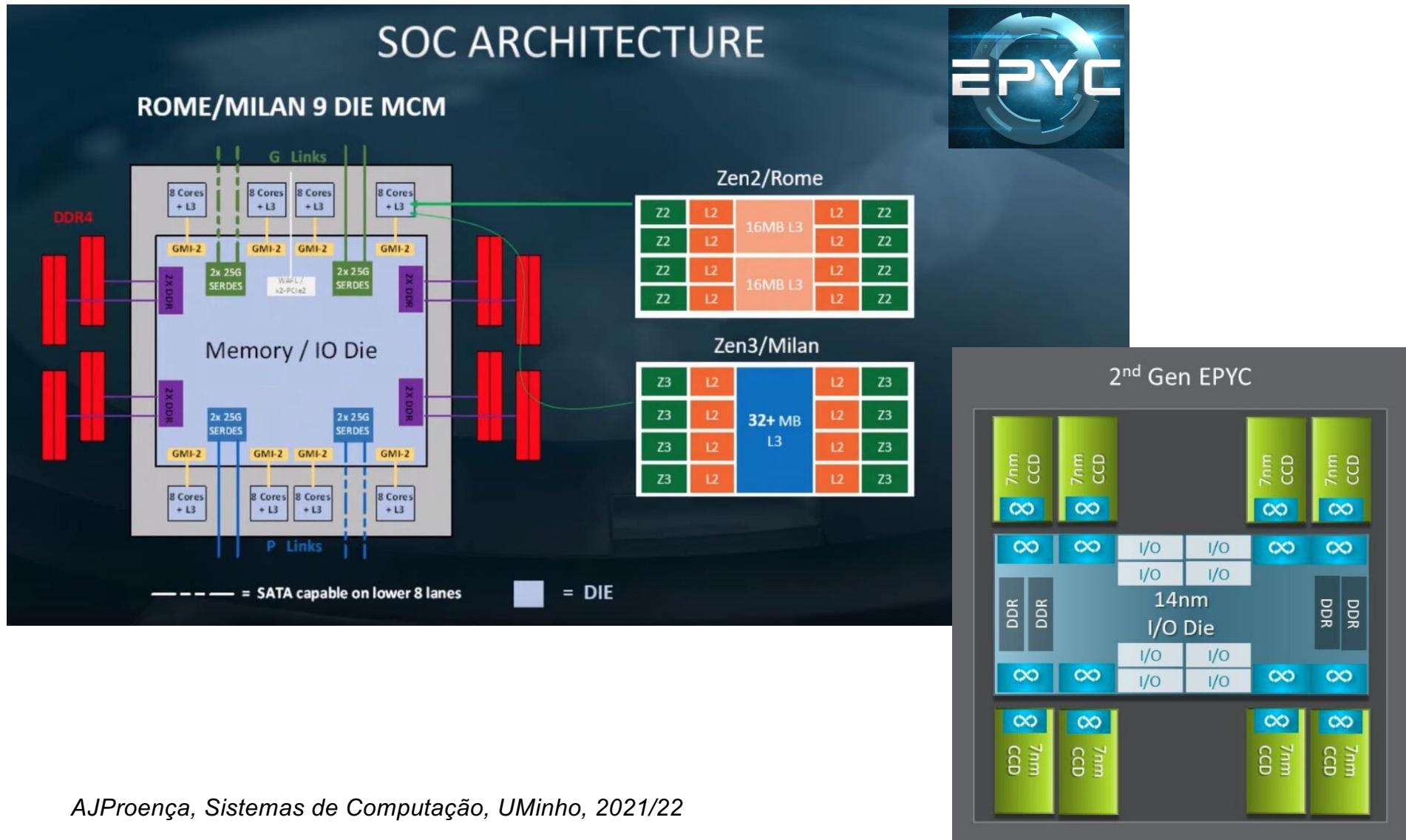
A concorrência da Intel nos servidores para HPC



- Em compatíveis com x86-64: AMD Epyc
 - processo de fabrico: 10nm (Intel) vs. 7nm (AMD)
 - máx nº de cores: 56 (Intel) vs. 64 (AMD)
- Em *chips* com design europeu: ARM V.8 (64-bits)
 - vários fabricantes com *dual-socket*
- Em *chips* chineses: Sunway, Huawei, ...
 - mais de 256 cores
 - coprocessadores com 128 cores
- Em coprocessadores japoneses: PEZY-2 (~2Ki cores)
- Recorde atual em #cores: Cerebras (400,000 cores)

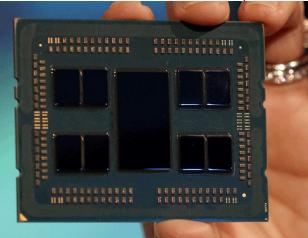


AMD Epyc: Zen 2 (Rome) & Zen 3 (Milan)

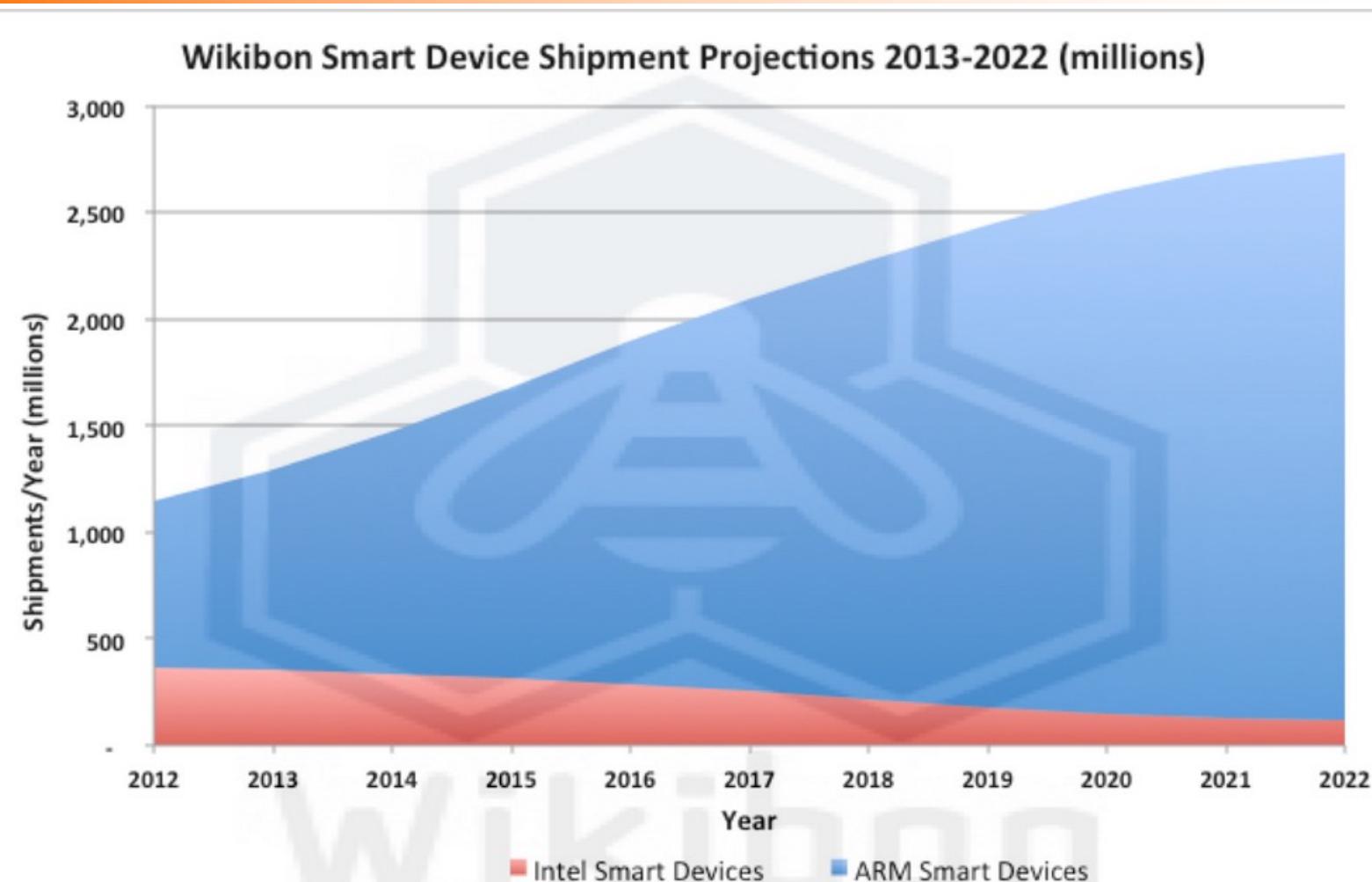


Intel Xeon vs. AMD Epyc

(ambos anunciados na 1^a metade 2019)

	AMD EPYC 7742 (Rome)	Intel Xeon Platinum 8280
7 nm, I/O 14 nm 9-die package		
2.25 GHz – 3.4 GHz		14 nm single-die, 2 SNC 2.7 GHz – 4.0 GHz AVX512: 1.8 – 3.7 GHz
Cores / Threads	64c / 128t	28c / 56t
L2/L3 Cache	512 KiB/core / 256 MB	1MiB/core / 38.5 MB
Max Memory/Bandwidth	4 TB / 190.7 GiB/s	1 TB / 131.13 GiB/s
Memory Channels	8	6
PCI-E Lanes	128x PCI-E 4.0	48x PCI-E 3.0
		 14 nm 2-die 82xx, 56 cores 12 memory channels

Processadores Intel vs. ARM



Source: Wikibon 2013, IDC & Gartner 2012 shipments & Wikibon 2013-2022 projections. Assumption: Apple & Microsoft migrate to successful 64-bit ARM.



WIKIPEDIA
The Free Encyclopedia



ARM brand: a bit of history...

ARM architecture

From Wikipedia, the free encyclopedia

ARM, previously **Advanced RISC Machine**, originally **Acorn RISC Machine**, is a family of reduced instruction set computing (RISC) architectures for computer processors, configured for various environments. Arm Holdings develops the architecture and licenses it to other companies, who design their own products that implement one of those architectures—including systems-on-chips (SoC) and systems-on-modules (SoM) that incorporate memory, interfaces, radios, etc. It also designs cores that implement this instruction set and licenses these designs to a number of companies that incorporate those core designs into their own products.

Processors that have a RISC architecture typically require fewer transistors than those with a complex instruction set computing (CISC) architecture (such as the x86 processors found in most personal computers), which

Current owner of Arm Holdings: NVidia

ARM architectures

The ARM logo, consisting of the lowercase word "arm" in a bold, blue, sans-serif font.

The ARM logo

Designer Arm Holdings

Bits 32-bit, 64-bit

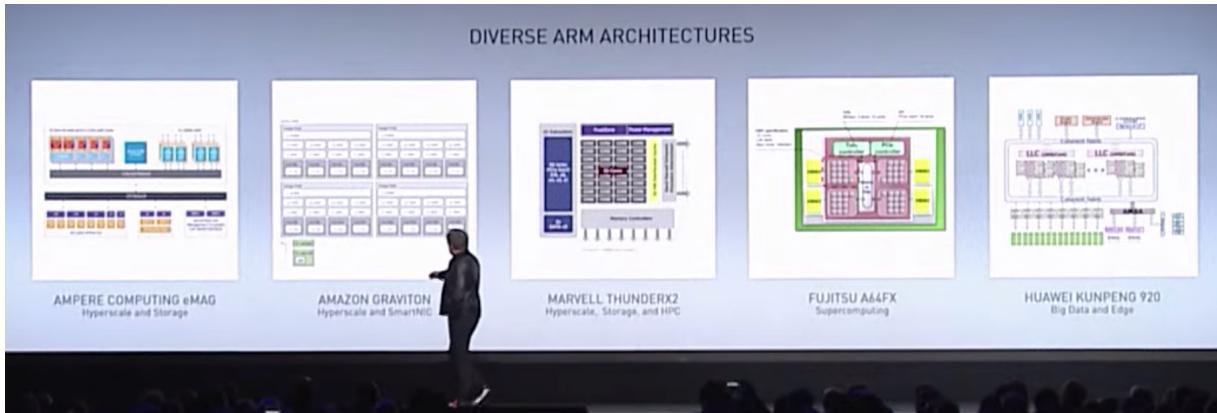
Introduced 1985; 34 years ago

Design RISC

Type Register-Register

Branching Condition code, compare and branch

Open Proprietary



HPCs with ARMv8: server-level competitors



FUJITSU



1. Marvell ThunderX product family

2. Fujitsu A64FX Arm chip

3. Ampere eMAG 8180 Arm Processor

4. Neoverse N1 hyperscale reference design

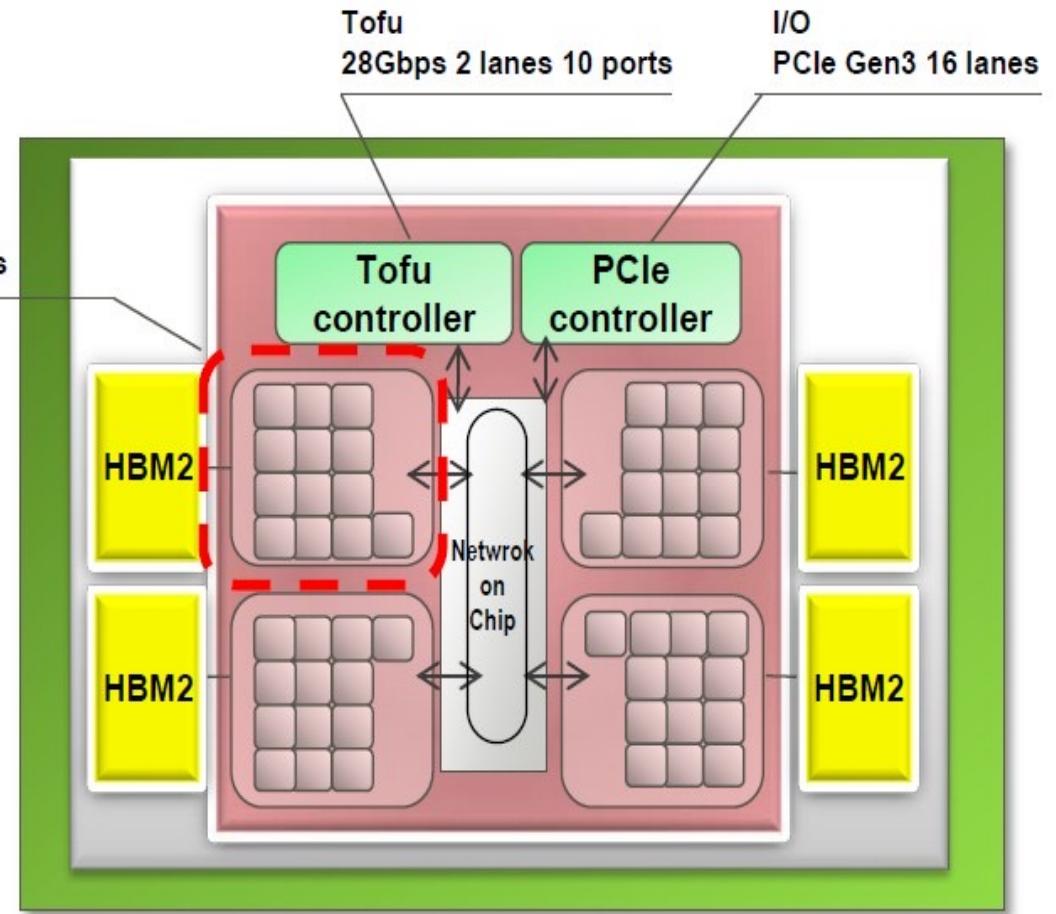
5. Huawei HiSilicon Kunpeng 920



Fujitsu's A64FX Arm Chip: 48+4 cores

A64FX Arm

- Feature size: 7 nm
 - Armv8.2-A spec with 512-bit SVE extensions
 - HP math and a dot-product engine
 - 4 core-memory groups
 - NoC: a double ring bus
 - cores in CMG linked by a crossbar to L2 cache & to HBM2 mem controller
 - 8 MiB L2 cache; no L3 cache
 - a Tofu-D controller on the die
 - #1 TOP500 since Jun'20 uses A64FX package
- CMG specification
13 cores
L2\$ 8MiB
Mem 8GiB, 256GB/s

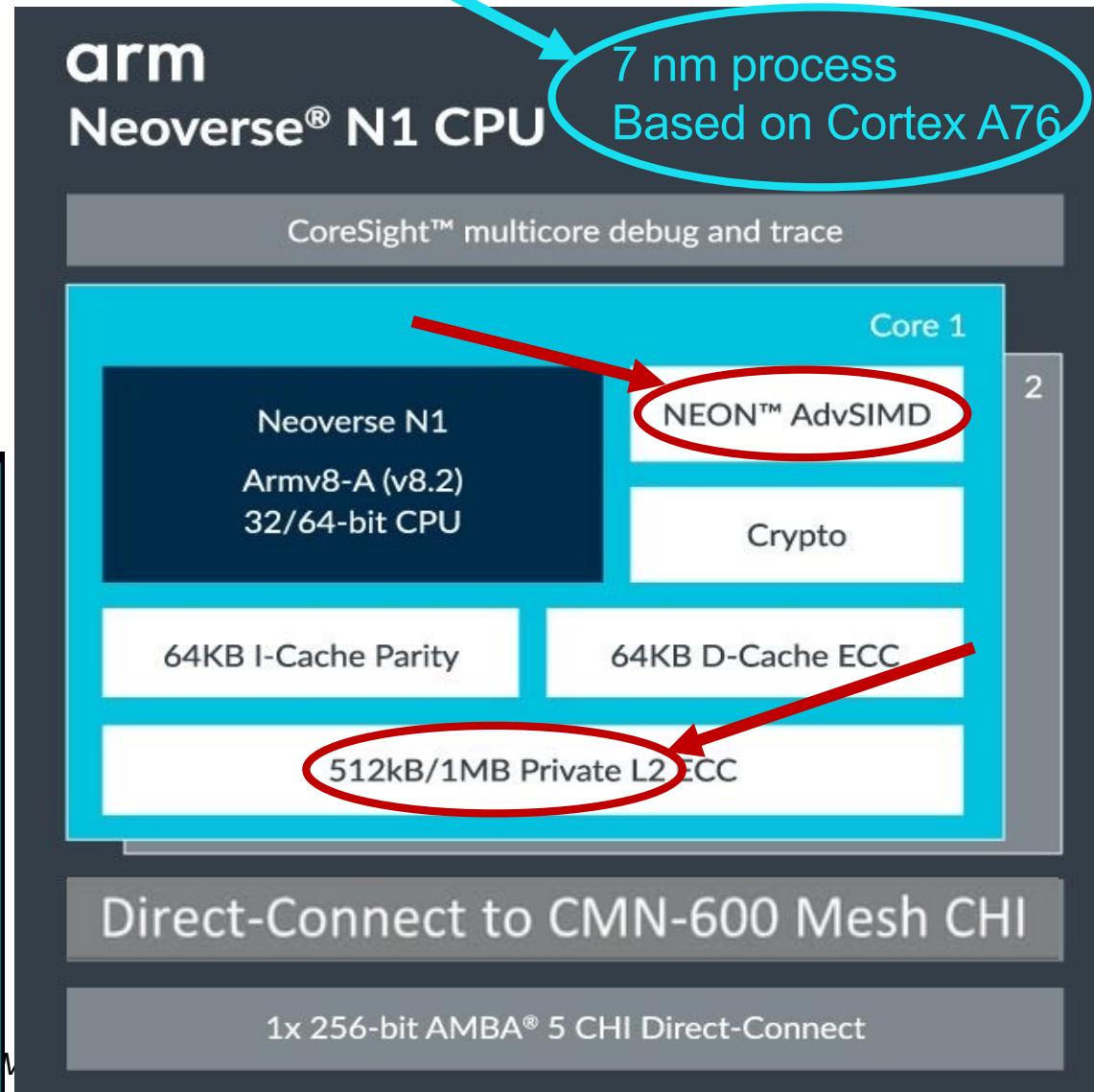
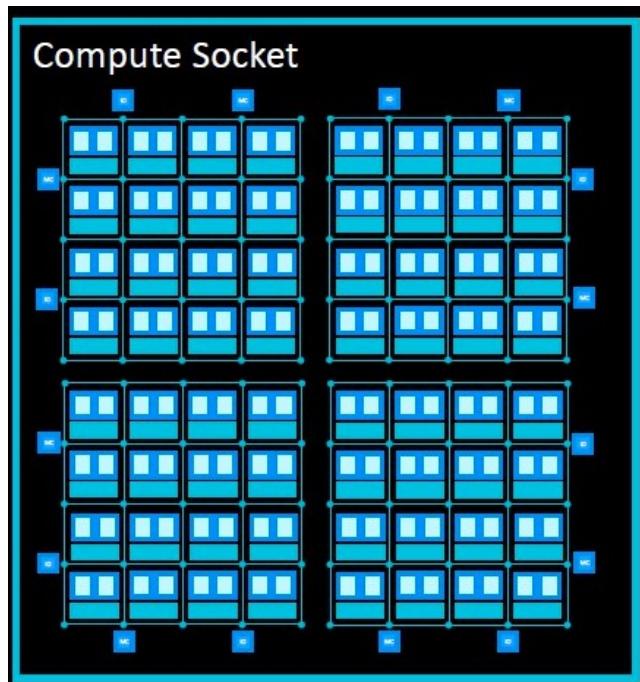
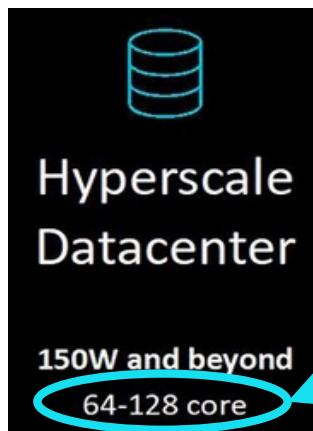




Arm Neoverse N1

Amazon (AWS) developed Graviton2 based on 64-core Neoverse

(announced Feb'19)





Sunway TaihuLight

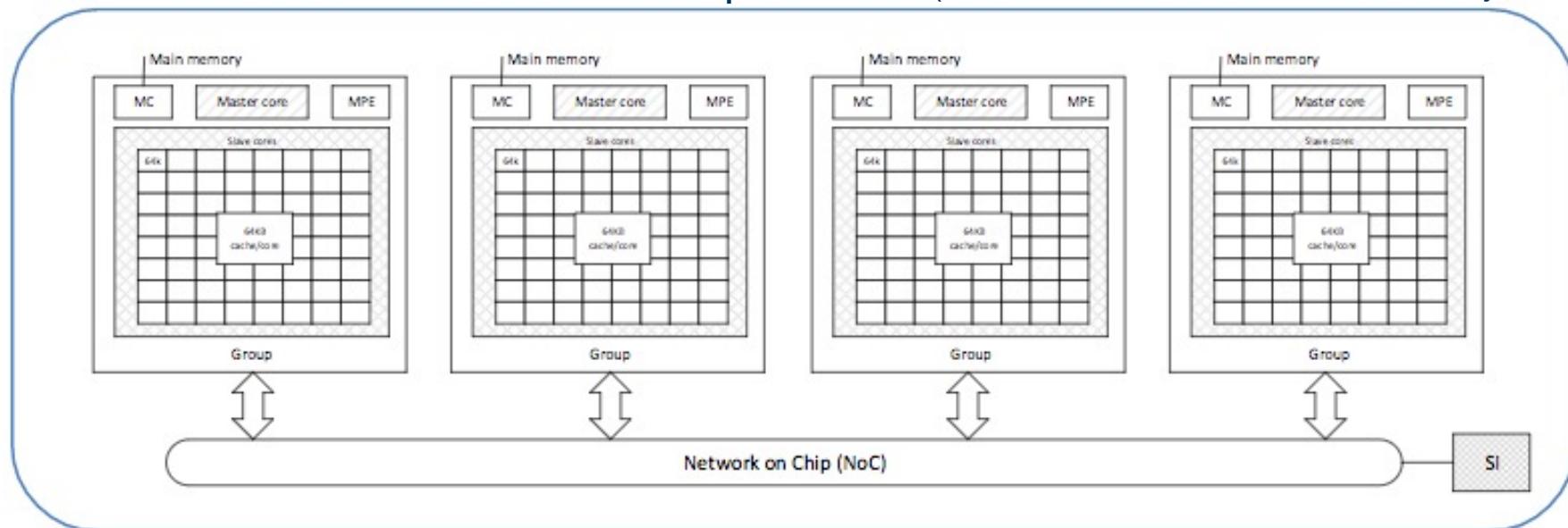
(#1 in June'16 TOP500)



One card with two nodes
(two SW26010 chips)



SW26010: the 4x64-core 64-bit RISC processor (w/ 256-bit vector instructions & only cache L1)



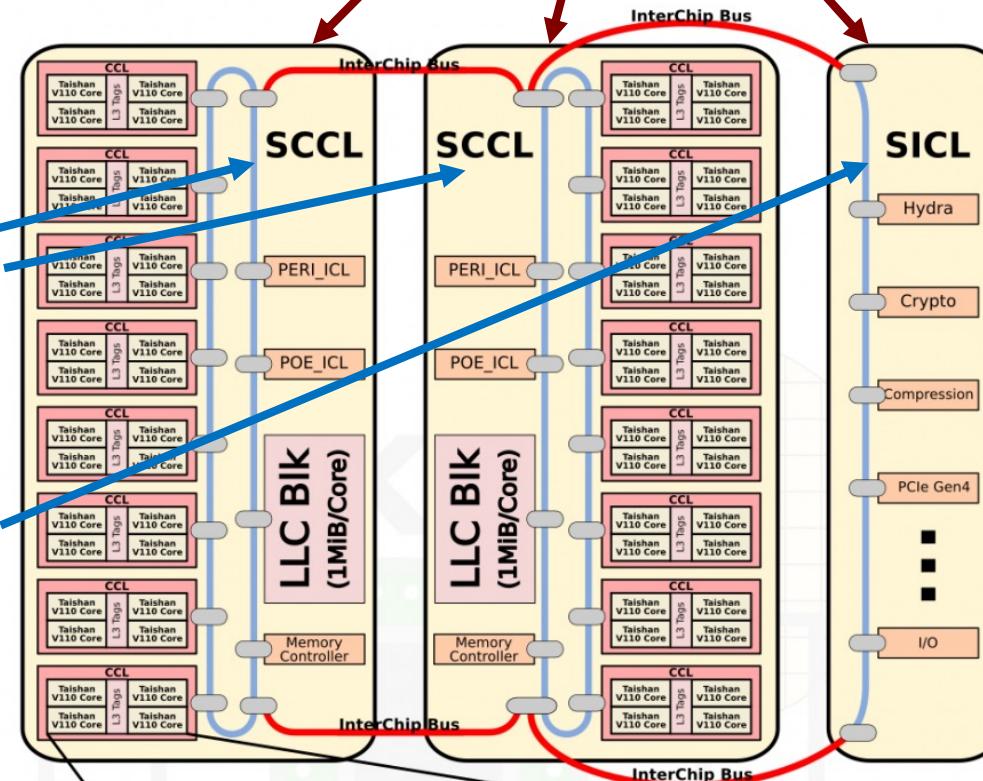


The Huawei Kunpeng 920: a multi-chip 48-64 cores

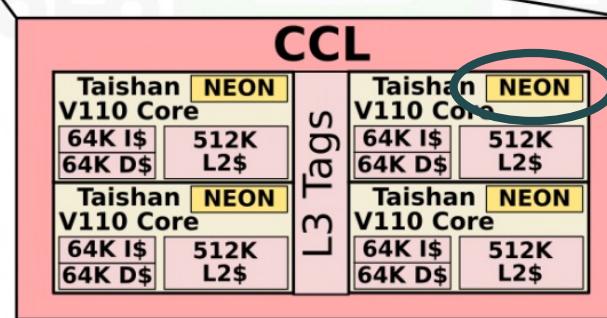
SCCL:
Super CPU Cluster

SICL:
Super IO Cluster

CCL:
CPU Clusters



128-bit
SIMD unit





Replacing the KNC in Tianhe-2A: the Matrix-2000 accelerator

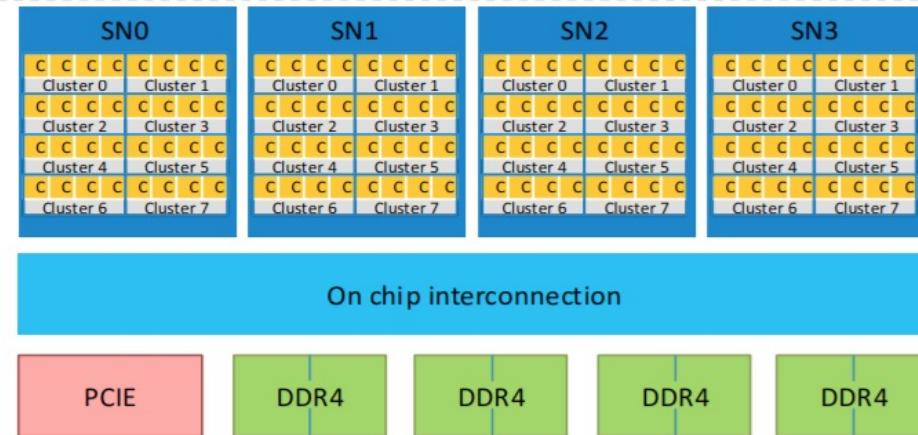


Matrix-2000 accelerator



Chip specification

- 128cores
 - 4 super-nodes (SN)
 - 8 clusters per SN
 - 4 cores per cluster
 - Core
 - Self-defined 256-bit vector ISA
 - 16 DP flops/cycle per core
- Peak performance: 2.4576Tflops@1.2GHz
 - 4 SNs x 8 clusters x 4cores x 16 flops x 1.2 GHz = 2.4576 Tflops



- Peak power dissipation: ~240W
- Interface
 - 8 DDR4-2400 channels
 - X16 PCIE 3.0 EP Port

10,000 PEZY-SC2 + 1,250 16-cores Xeon =
19.84 M PEZY cores + 20 K Xeon cores



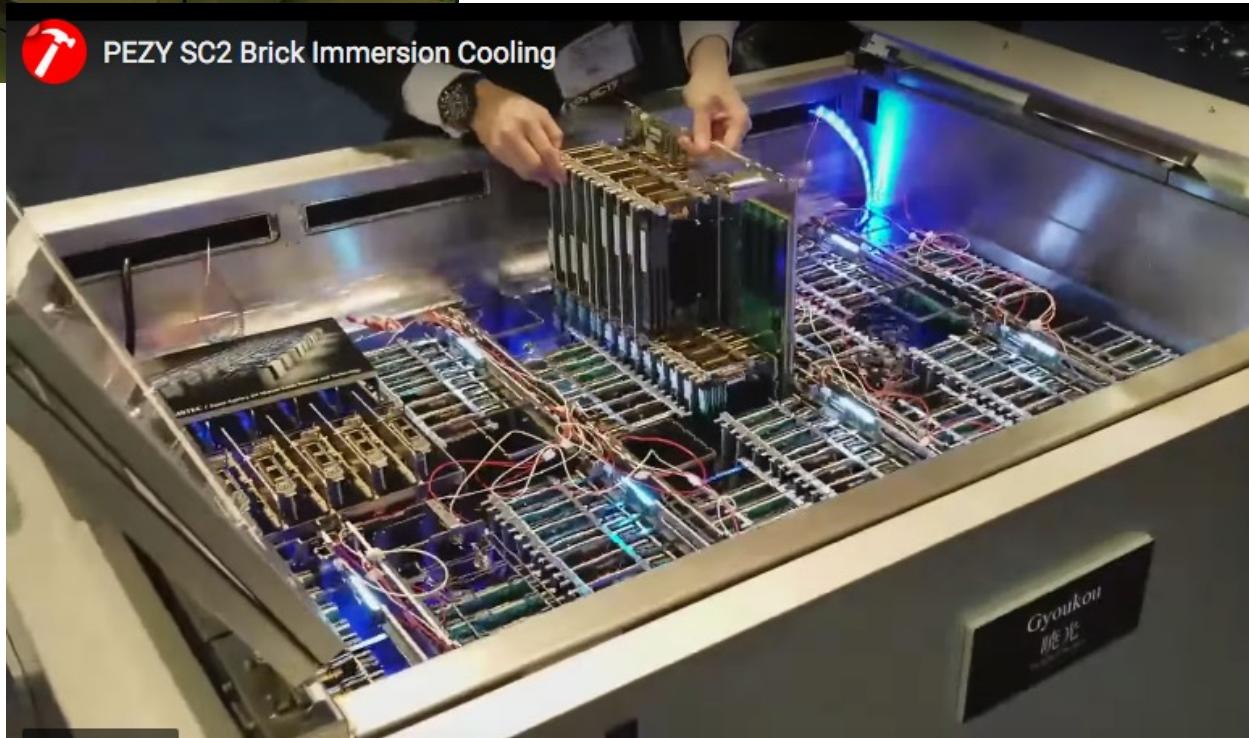
Gyoukou ZettaScaler-2.2

Nov'17 Green500

5	4	Gyoukou - ZettaScaler-2.2	19,860,000	19.1
		HPC system, Xeon D-1571		
		16C 1.3GHz, Infiniband		
		EDR, PEZY-SC2 700Mhz ,		
		ExaScaler		



20 immersion tanks
each tank 16 bricks
each brick 32 PEZY
each PEZY ~2K
8-way SMT cores
=>
each tank ~1M cores





Cerebras Wafer Scale Engine (WSE): the largest chip ever built



46,225 mm² chip

56x larger than the biggest GPU ever made

400,000 core

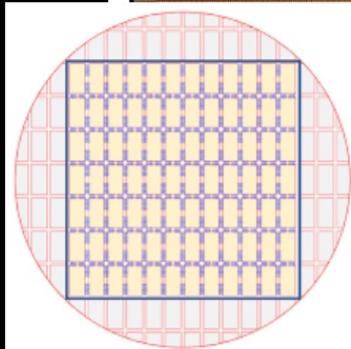
78x more cores

18 GB on-chip SRAM

3000x more on-chip memory

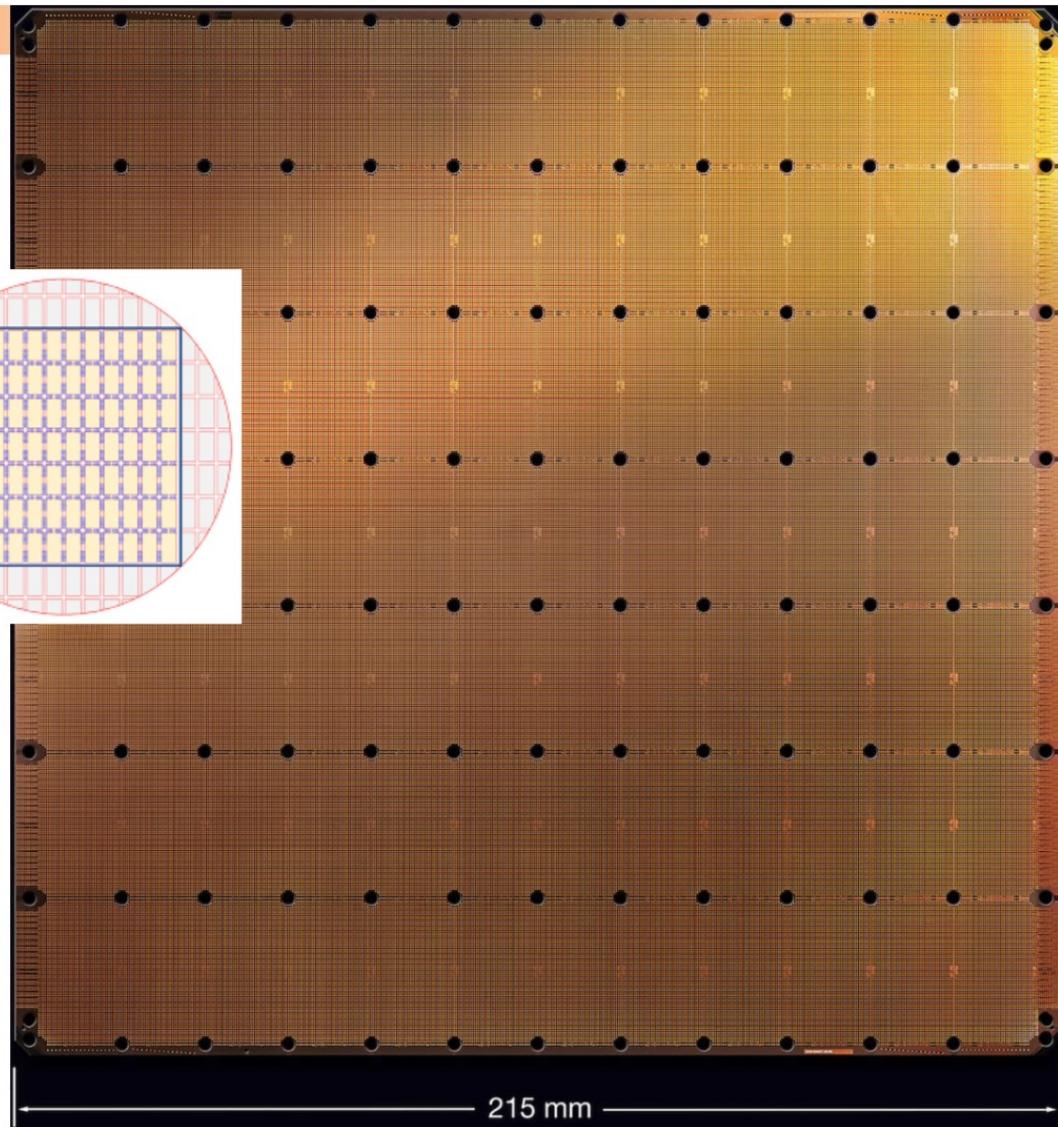
100 Pb/s interconnect

33,000x more bandwidth



**400,000
Sparse Linear Algebra
(SLA) cores
(AI-optimized)**

AJProen , Advanced Architectures, MiEI,



Aceleradores de computação: tipos de componentes (devices)

Hardware acceleration

From Wikipedia, the free encyclopedia

(Redirected from [Hardware accelerator](#))

In [computing](#), **hardware acceleration** is the use of [computer hardware](#) specially made to perform some functions more efficiently than is possible in [software](#) running on a general-purpose [CPU](#).

Tipos mais comuns de componentes para acelerar funções específicas:

- **DSP**: *Digital Signal Processor*, para processar sinais elétricos 1D, codificados em binário, por ex. para tlm, TV, ...
- **GPU**: *Graphics Processing Unit*, também usado em computação científica (sem interfaces gráficas), com milhares de unidades de FP
- **NNP**: *Neural Net Processor*, usado em aplicações de AI, nomeadamente para treino e classificação com redes neurais
- **FPGA**: *Field Programmable Gate Arrays*, blocos de unidades lógicas rudimentares que permitem construir qq circuito para executar uma operação