Master Informatics Eng.

2018/19 *A.J.Proença*

Data Parallelism 2 (SIMD++, NVidia GPUs...) (most slides are borrowed)

AJProença, Advanced Architectures, MiEI, UMinho, 2018/19

Beyond Vector/SIMD architectures

XX

XX

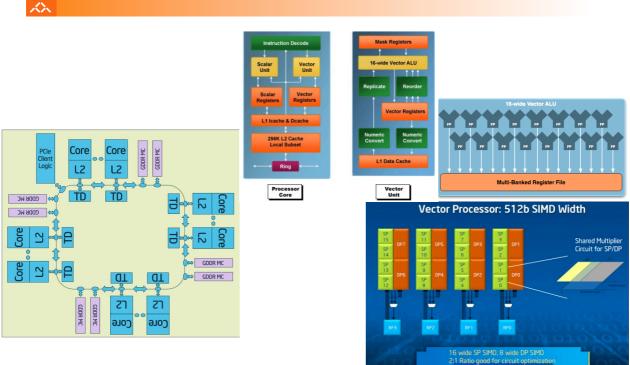
- Vector/SIMD-extended architectures are hybrid approaches
 - mix (super)scalar + vector op capabilities on a single device
 - highly pipelined approach to reduce memory access penalty
 - tightly-closed access to shared memory: lower latency
- Evolution of Vector/SIMD-extended architectures
 - PU (Processing Unit) cores with wider vector units
 - <u>x86</u> many-core: Intel MIC / Xeon KNL
 - ...
 - coprocessors (require a host scalar processor): accelerator devices
 - on disjoint physical memories (e.g., Xeon KNC with PCI-Expr, ...)
 - ...
 - heterogeneous PUs in a SoC: multicore PUs with GPU-cores
 - ...

1

Intel MIC: Many Integrated Core

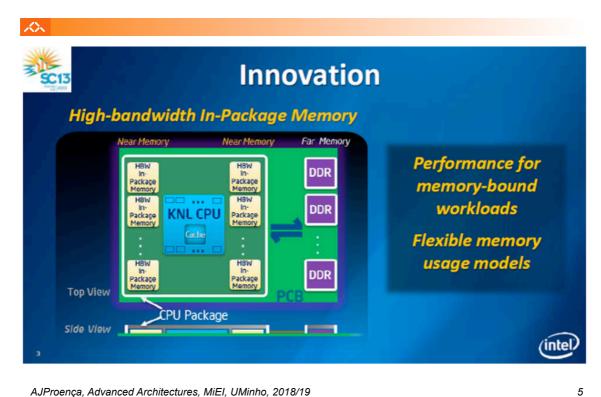


Intel Knights Corner architecture



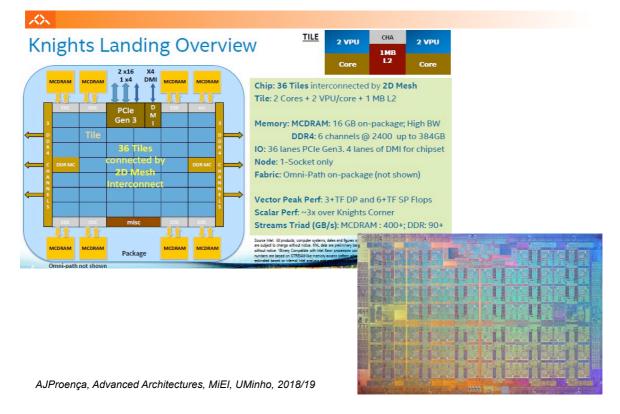
AJProença, Advanced Architectures, MiEl, UMinho, 2018/19

The new Knights Landing architecture



AJProença, Advanced Architectures, MiEI, UMinho, 2018/19

Intel Knights Landing in 2016: Xeon Phi com 72 active cores

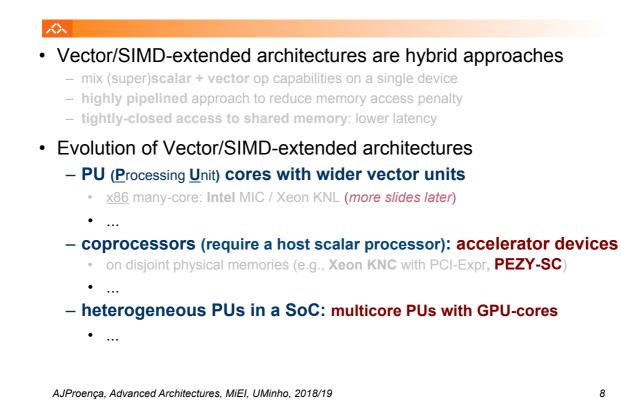


INTEL[®] XEON PHI[™] X200 PROCESSOR OVERVIEW





Beyond Vector/SIMD architectures



PEZY-SC: <u>P</u>eta_<u>Exa_Zetta_Y</u>otta-<u>SuperComputer:</u> a 1024-core many-core processor chip

Green500 Rank	MFLOPS/W	Site*	Computer*		Green500 list	
1 94	6,673.84	Advanced Center for Computing and Communication, RIKEN	Shoubu - ZettaSener-1.6, Xoon E5-2618Lv3 8C 2.3GHz, Infiniband FDR PEZY-SCnp	149.99	June'2016	
² 486	6,195.22	Computational Astrophysics Laboratory, RIKEN	Satsuki - ZettaSca ler - 1.6, Ks on E5-2618Lv3 8C 2.3GHz, Infiniband FDR PEZY-SCnp	46.89		
3 ₁	6,051.30	National Supercomputing Center in Wuxi	Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway	15,371.00		
4 440	5,272.09	GSI Helmholtz Center	ASUS ESC4000 FDR/G2S, Intel Xeon E5-2690v2 10C 3GHz Infiniband FDR, AMD FirePro S9150	57.15 City		
⁵ 446	4,778.46	Institute of Modern Physics (IMP), Chinese Academy of Sciences	Sugon Cluster W780I, Xeon E5-2640v3 8C 2.6GH QDR, NVIDIA Tesla K80	(16PE)	Name in the second second	
6 122	4,112.11	Stanford Research Computing Center	Sugon Cluster W7801, Xeon E5-2640v3 8C 2.6GF QDR, NVIDIA Tesla K80 XStream - Cray CS-Storm, Intel Xeon E5-2680v2 Infiniband FDR, Nvidia K80			
Top500 Rank			PETY-SC	Pre	23 \$	
			PEZY-SC And Generation Many Core Processor with 1024. Cores Supported by 2013 NEDD Project PEZY Computing K.K. B27701432-ES			
AJF	Proença, .	Advanced Architectures, N	ліЕІ, UMinho, 2018/19	PC	IEGen3 ARM Prefecture	

Beyond Vector/SIMD architectures

1

Vector/SIMD-extended architectures are hybrid approaches

- mix (super)scalar + vector op capabilities on a single device
- highly pipelined approach to reduce memory access penalty
- tightly-closed access to shared memory: lower latency

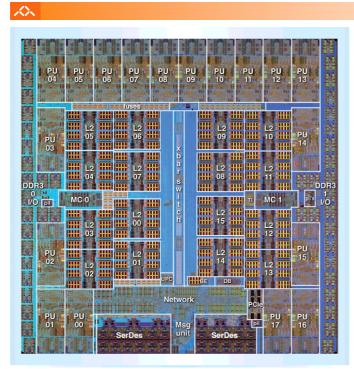
Evolution of Vector/SIMD-extended architectures

- PU (Processing Unit) cores with wider vector units
 - <u>x86</u> many-core: Intel MIC / Xeon KNL (more slides later)
 - other many-core: IBM Power BlueGene/Q Compute, ShenWay 260
- coprocessors (require a host scalar processor): accelerator devices
 - on disjoint physical memories (e.g., Xeon KNC with PCI-Expr, PEZY-SC)
 - .

- heterogeneous PUs in a SoC: multicore PUs with GPU-cores

• ...

IBM Power BlueGene/Q Compute (chip)

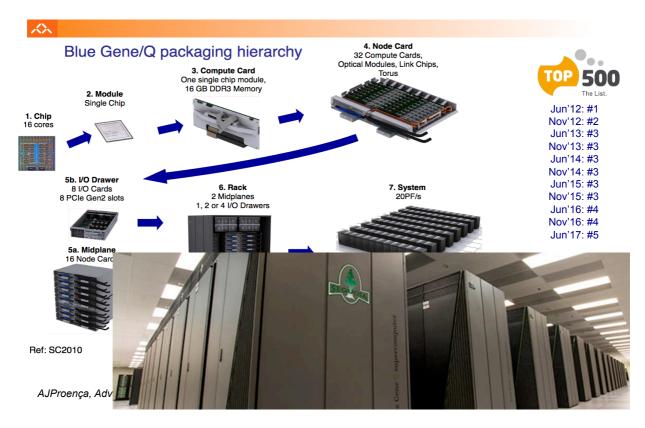


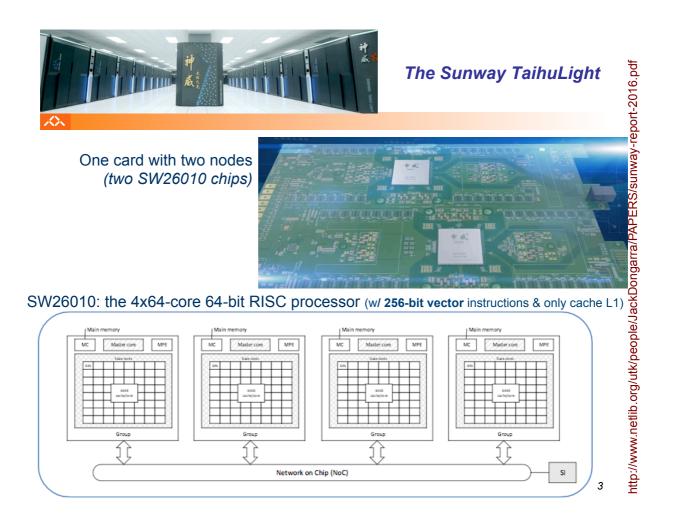
AJProença, Advanced Architectures, MiEl, UMinho, 2018/19

Features:

- launched in 2010/11 (TOP500: #1 in Jun12, #4 in Jun16)
- 18-cores
 - 16 compute,
 1 OS support, 1 redundant
 - 64 bits PowerISA
 - 1.6 GHz
 - L1 I/D cache => 16 kB / 16 kB
 - each core: <u>quad-FPU</u> (4-wide double precision SIMD)
 - each core: 4-way multi-threaded
- shared L2 cache: 32 MB
- dual memory controller
- IBM ended development of BlueGene project in 2015...

IBM Power BlueGene/Q Compute (Sequoia system)





Beyond Vector/SIMD architectures

\sim

Vector/SIMD-extended architectures are hybrid approaches

- mix (super)scalar + vector op capabilities on a single device
- highly pipelined approach to reduce memory access penalty
- tightly-closed access to shared memory: lower latency

Evolution of Vector/SIMD-extended architectures

- PU (Processing Unit) cores with wider vector units

- <u>x86</u> many-core: Intel MIC / Xeon KNL (more slides later)
- other many-core: IBM BlueGene/Q Compute, ShenWay 260
- coprocessors (require a host scalar processor): accelerator devices
 - on disjoint physical memories (e.g., Xeon KNC with PCI-Expr, PEZY-SC)
 - ISA-free architectures, code compiled to silica: FPGA
 - ...

– heterogeneous PUs in a SoC: multicore PUs with GPU-cores

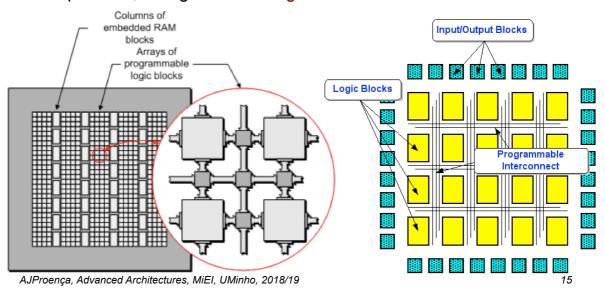
• ...

What is an FPGA

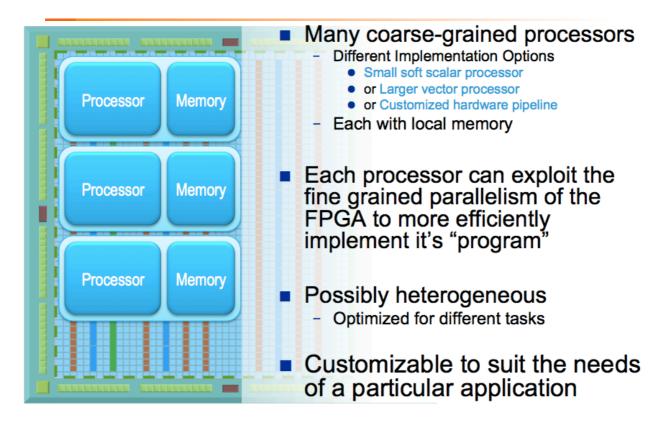
\sim

Field-Programmable Gate Arrays (FPGA)

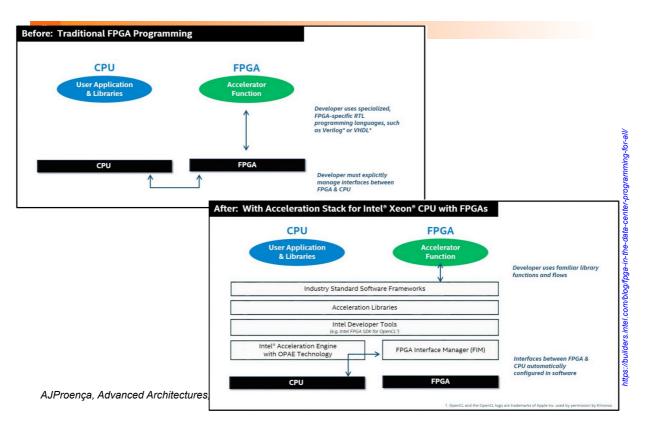
A fabric with 1000s of simple configurable logic cells with LUTs, on-chip SRAM, configurable routing and I/O cells



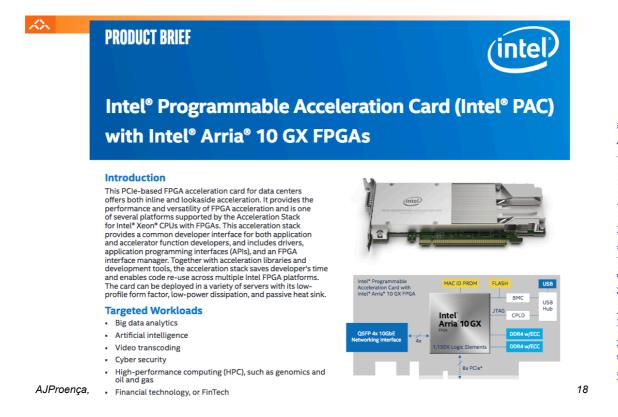
FPGA as a multiple configurable ISA



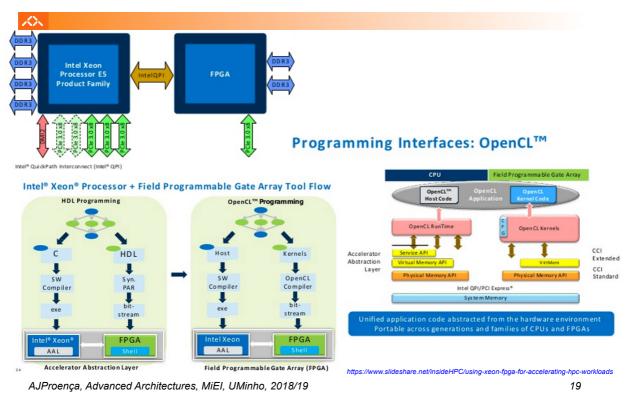
FPGA as a computing accelerator



The Intel Programmable Acceleration Card



Faster integration of programmable acceleration cards at Intel



Beyond Vector/SIMD architectures

~~

Vector/SIMD-extended architectures are hybrid approaches

- mix (super)scalar + vector op capabilities on a single device
- highly pipelined approach to reduce memory access penalty
- tightly-closed access to shared memory: lower latency

Evolution of Vector/SIMD-extended architectures

- PU (Processing Unit) cores with wider vector units

- <u>x86</u> many-core: Intel MIC / Xeon KNL (more slides later)
- other many-core: IBM BlueGene/Q Compute, ShenWay 260

- coprocessors (require a host scalar processor): accelerator devices

- on disjoint physical memories (e.g., **Xeon KNC** with PCI-Expr, **PEZY-SC**)
- ISA-free architectures, code compiled to silica: FPGA
- focus on SIMT/SIMD to hide memory latency: GPU-type approach
- ..

٠

- heterogeneous PUs in a SoC: multicore PUs with GPU-cores

AJProença, Advanced Architectures, MiEl, UMinho, 2018/19

Graphical Processing Units

Graphical Processing Units

- Question to GPU architects:
 - Given the hardware invested to do graphics well, how can we supplement it to improve the performance of a wider range of applications?
- Key ideas:
 - Heterogeneous execution model
 - CPU is the host, GPU is the device
 - Develop a C-like programming language for GPU
 - Unify all forms of GPU parallelism as CUDA threads
 - Programming model follows SIMT: "Single Instruction Multiple Thread"

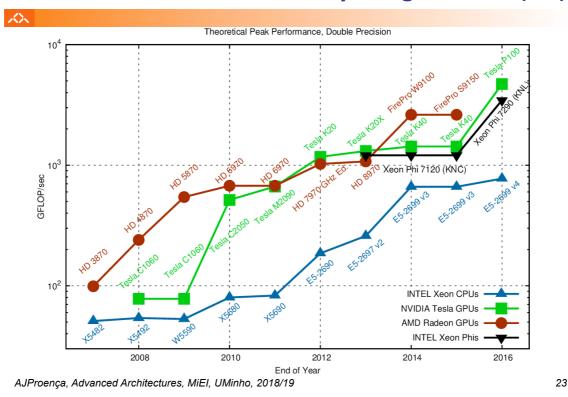
Copyright © 2012, Elsevier Inc. All rights reserved.

21

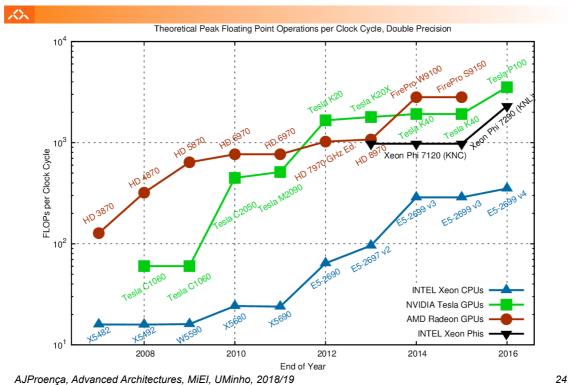
http://www.karlrupp.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time. # cores/processing elements in several devices 1 Number of Physical Cores/Multiprocessors, High-End Hardware 10² Key question: Xeon Phi 7290 (KNL what is a core? Xeon Phi 7120 (KNC GTX 28 HD 6970 HD 7970 a) IU+FPU? GPU-type... ન્સ્ઈ GTX 58 10¹ E5-2697 4870 E5-2690 HD 3870 ×5690 and in this course: 15492 W5590 - b) INTEL Xeon CPUs NVIDIA GeForce GPUs Note: the web link AMD Radeon GPUs with these plots was **INTEL Xeon Phis** 10⁰ updated in Aug'16 2008 2010 2012 2014 2016 End of Year 22

AJProença, Advanced Architectures, MiEI, UMinho, 2018/19

Theoretical peak performance in several computing devices (DP)



Theoretical peak FP Op's per clock cycle in several computing devices (DP)



http://www.karlrupp.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time/

NVIDIA GPU Architecture

Similarities to vector machines:

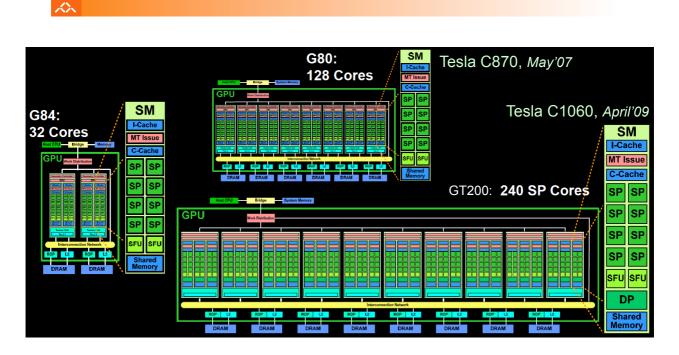
- Works well with data-level parallel problems
- Scatter-gather transfers
- Mask registers
- Large register files

Differences:

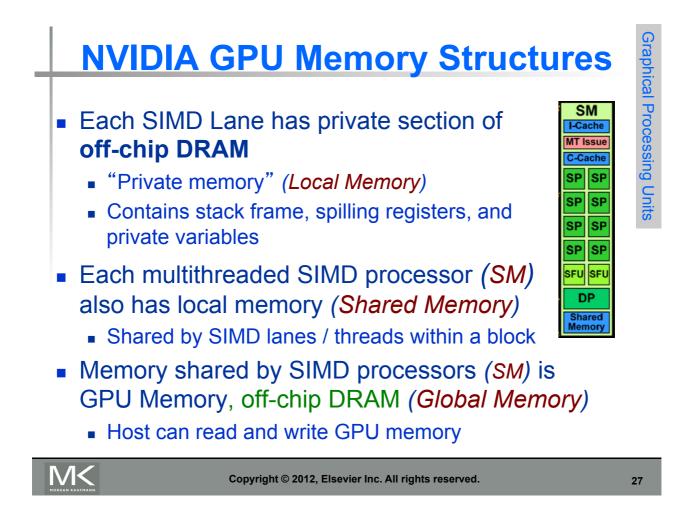
- No scalar processor
- Uses multithreading to hide memory latency
- Has many functional units, as opposed to a few deeply pipelined units like a vector processor



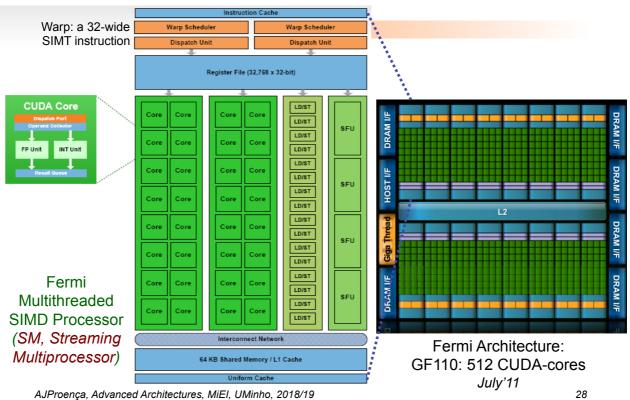
Early NVidia GPU Computing Modules



AJProença, Advanced Architectures, MiEl, UMinho, 2018/19



The NVidia Fermi architecture



Fermi Architecture Innovations

Each SIMD processor has

- Two SIMD thread schedulers, two instruction dispatch units
- 16 SIMD lanes (SIMD width=32, chime=2 cycles), 16 load-store units, 4 special function units
- Thus, two threads of SIMD instructions are scheduled every two clock cycles

Instruction Cache				
Warp Scheduler	Warp Scheduler			
Dispatch Unit	Dispatch Unit			
+	+			

- Fast double precision
- Caches for GPU memory (16/64KB_L1/SM and global 768KB_L2)
- 64-bit addressing and unified address space
- Error correcting codes
- Faster context switching
- Faster atomic instructions

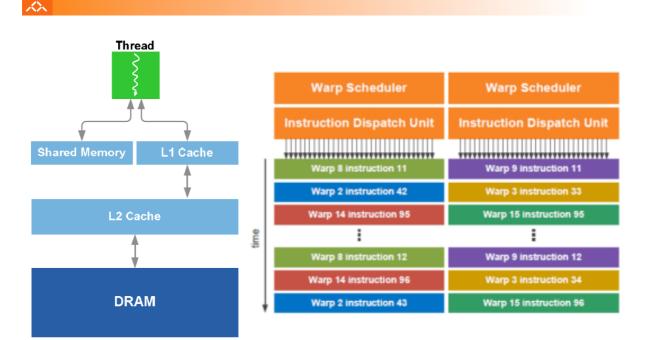
AJProença, Advanced Architectures, MiEI, UMinho, 2018/19

Copyright © 2012, Elsevier Inc. All rights reserved.

29

Graphical Processing Units

Fermi: Multithreading and Memory Hierarchy



TOP500 list in November 2010: 3 systems in the top4 use Fermi GPUs

 \mathcal{K}



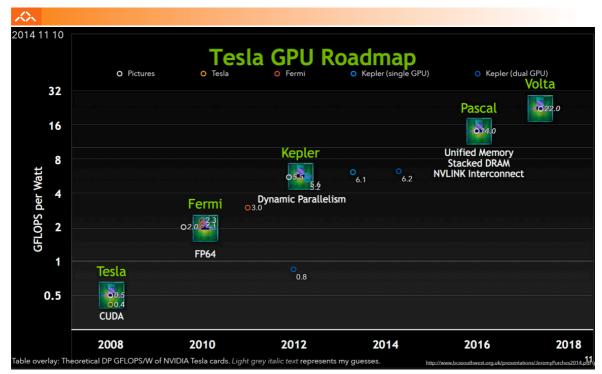
HIGHLIGHTS: NOVEMBER 2010

- The Chinese Tianhe-1A system is the new No. 1 on the TOP500 and clearly in the lead with 2.57 petaflop/s performance.
- No. 3 is also a Chinese system called Nebulae, built from a Dawning TC3600 Blade system with Intel X5650 processors and NVIDIA Tesla C2050 GPUs
- There are seven petaflop/s systems in the TOP10
- The U.S. is tops in petaflop/s with three systems performing at the petaflop/s level
- The two Chinese systems and the new Japanese Tsubame 2.0 system at No. 4 are all using NVIDIA GPUs to accelerate computation and a total of 28 systems on the list are using GPU technology.

AJProença, Advanced Architectures, MiEI, UMinho, 2018/19

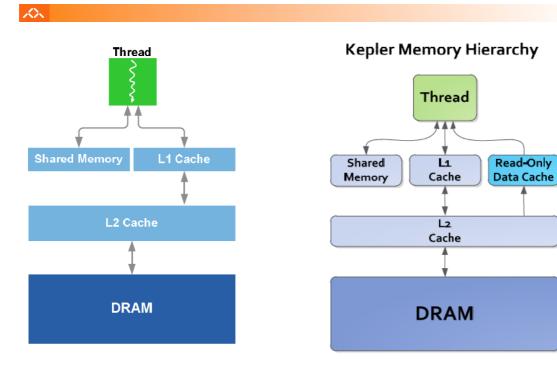
31

Families in NVidia Tesla GPUs

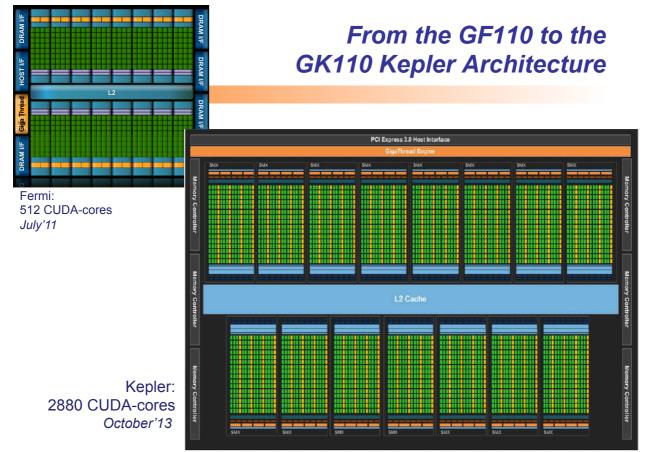


AJProença, Advanced Architectures, MiEI, UMinho, 2018/19

From Fermi into Kepler: The Memory Hierarchy



AJProença, Advanced Architectures, MiEI, UMinho, 2018/19



AJProença, Advanced Architectures, MiEI, UMinho, 2018/19

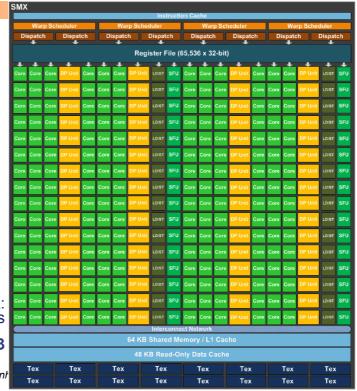
33

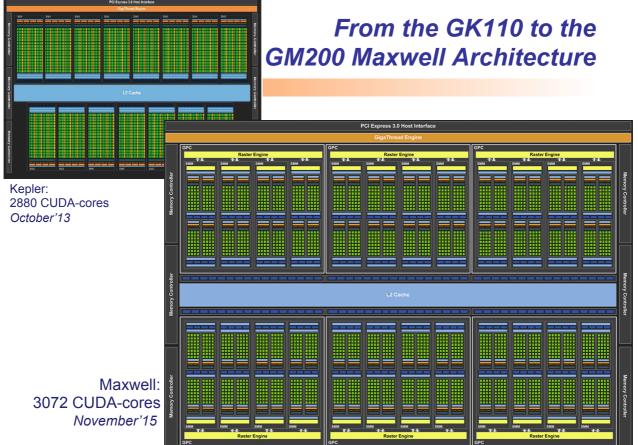


192 CUDA-cores Ratio DPunit : SPunit -> 1 : 3

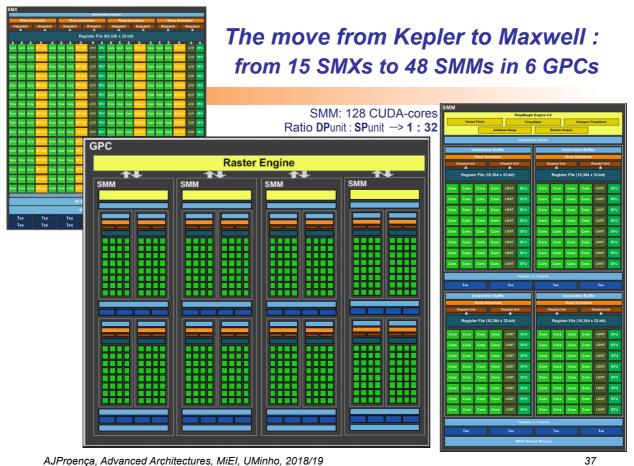
AJProença, Advanced Architectures, MiEl, UMinł

From Fermi to Kepler core: SM and the SMX Architecture

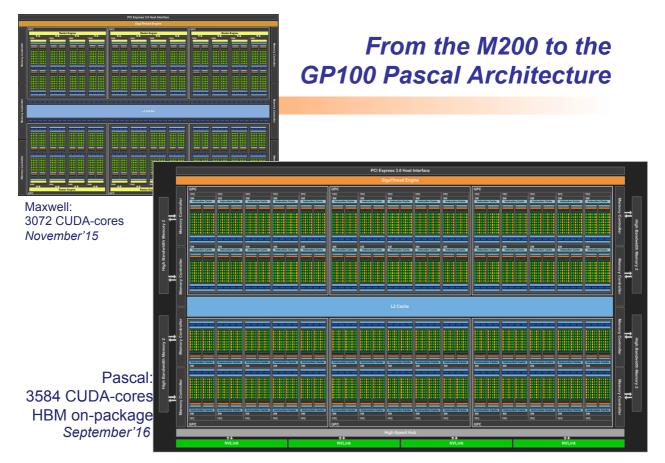




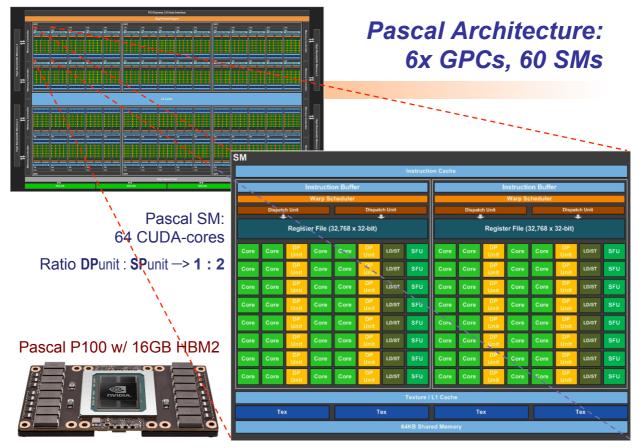
AJProença, Advanced Architectures, MiEI, UMinho, 2018/19



AJProença, Advanced Architectures, MiEI, UMinho, 2018/19



AJProença, Advanced Architectures, MiEl, UMinho, 2018/19



AJProença, Advanced Architectures, MiEI, UMinho, 2018/19

<section-header>Pasca: Status CuDA-cores November 15

AJProença, Advanced Architectures, MiEl, UMinho, 2018/19

39

	Volta Architecture: -6x-GPCs, 80 SMs					
	L0 Instruction Cache	L0 Instruction Cache				
	Warp Scheduler (32 thread/clk)	Warp Scheduler (32 thread/clk)				
	Dispatch Unit (32 thread/clk) Register File (16,384 x 32-bit)	Dispatch Unit (32 thread/clk) Register File (16,384 x 32-bit)				
	FP64 INT INT FP32 FP32 FP64 INT INT FP32 FP32	FP64 INT INT FP32 FP32 FP64 INT INT FP32 FP32				
Hiphford Huk MC Hik MC Hik	FP64 INT INT FP32 FP32	FP64 INT INT FP32 FP32				
	FP64 INT INT FP32 FP32 TENSOR TENSOR	FP64 INT INT FP32 FP32 TENSOR TENSOR				
Volta SM:	FP64 INT INT FP32 FP32 CORE CORE	FP64 INT INT FP32 FP32 CORE CORE				
64 CUDA-cores	FR64 INT INT FP32 FP32 FP64 INT INT FP32 FP32	FP64 INT INT FP32 FP32 FP64 INT INT FP32 FP32				
•	FP64 INT INT FP32 FP32	FP64 INT INT FP32 FP32				
New: 8 Tensor-cores	LD/ LD/ LD/ LD/ LD/ LD/ LD/ LD/ ST	LD/ LD/ LD/ LD/ LD/ LD/ LD/ LD/ ST				
Ratio DPunit : SPunit -> 1 : 2	L0 Instruction Cache	L0 Instruction Cache				
	Warp Scheduler (32 thread/clk) Dispatch Unit (32 thread/clk)	Warp Scheduler (32 thread/clk) Dispatch Unit (32 thread/clk)				
	Register File (16,384 x 32-bit)	Register File (16,384 x 32-bit)				
	FP64 INT INT FP32 FP32	FP64 INT INT FP32 FP32				
Volta V100 w/ 16GB HBM2	FP64 INT INT FP32 FP32	FP64 INT INT FP32 FP32				
	FP64 INT INT FP32 FP32	FP64 INT INT FP32 FP32				
	FP64 INT INT FP32 FP32 TENSOR TENSOR FP64 INT INT FP32 FP32 CORE CORE	FP64 INT INT FP32 FP32 TENSOR TENSOR FP64 INT INT FP32 FP32 CORE CORE				
	FP64 INT INT FP32 FP32	FP64 INT INT FP32 FP32				
	FP64 INT INT FP32 FP32	FP64 INT INT FP32 FP32				
	FP64 INT INT FP32 FP32	FP64 INT INT FP32 FP32				
	LD/ LD/ LD/ LD/ LD/ LD/ LD/ SFU	LD/ LD/ LD/ LD/ LD/ LD/ LD/ ST				
AJProenca, Advanced Architectures, MiEI, UM	128KB L1 Data Cache	e / Shared Memory				
··· ··· ······························	Iex Iex	Tex Tex				

Tesla Product	Tesla K40	Tesla M40	Tesla P100	Tesla V100
GPU	GK180 (Kepler)	GM200 (Maxwell)	GP100 (Pascal)	GV100 (Volta)
SMs	15	24	56	80
TPCs	15	24	28	40
FP32 Cores / SM	192	128	64	64
FP32 Cores / GPU	2880	3072	3584	5120
FP64 Cores / SM	64	4	32	32
FP64 Cores / GPU	960	96	1792	2560
Tensor Cores / SM	NA	NA	NA	
Tensor Cores / GPU	NA	NA	NA	640
GPU Boost Clock	810/875 MHz	1114 MHz	1480 MHz	1530 MHz
Peak FP32 TFLOP/s	5.04	6.8	10.6	15.7
Peak FP64 TFLOP/s	1.68	.21	5.3	7.8
Peak Tensor Core TFLOP/s	NA	NA	NA	125
Texture Units	240	1 92	224	320
Memory Interface	384-bit GDDR5	384-bit GDDR5	4096-bit HBM2	4096-bit HBM2
Memory Size	Up to 12 GB	Up to 24 GB	16 GB	16 GB
L2 Cache Size	1536 KB	3072 KB	4096 KB	6144 KB
Shared Memory Size / SM	16 KB/32 KB/48 KB	96 KB	64 KB	Configurable up to 96 KB
Register File Size / SM	256 KB	256 KB	256 KB	256KB
Register File Size / GPU	3840 KB	6144 KB	14336 KB	20480 KB
TDP	235 Watts	250 Watts	300 Watts	300 Watts
Transistors	7.1 billion	8 billion	15.3 billion	21.1 billion
GPU Die Size	551 mm²	601 mm²	610 mm ²	815 mm²
Manufacturing Process	28 nm	28 nm	16 nm FinFET+	12 nm FFN

Tesla accelerators: recent evolution

ANNOUNCING TESLA V100 GIANT LEAP FOR AI & HPC VOLTA WITH NEW TENSOR CORE 21B xtors | TSMC 12nm FFN | 815mm² 5,120 CUDA cores 7.5 FP64 TFLOPS | 15 FP32 TFLOPS NEW 120 Tensor TFLOPS 20MB SM RF | 16MB Cache | 16GB 4B 300 GB/s NVLink

https://devblogs.nvidia.com/parallelforall/inside-volta/