**Advanced Architectures** 

# Master Informatics Eng.

2017/18 *A.J.Proença* 

Data Parallelism 2 (SIMD++, Intel MIC, NVidia GPU ...) (most slides are borrowed)

AJProença, Advanced Architectures, MiEI, UMinho, 2017/18

XX

# **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

• ...

- 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

# Intel MIC: <u>Many Integrated Core</u>

### 公

# Intel evolution, from:

• Larrabee (80-core GPU)



# & SCC

<u>S</u>ingle-chip <u>C</u>loud <u>C</u>omputer, 24x dual-core tiles



# to MIC:

- Knights Ferry (pre-production, Stampede)
- Knights Corner Xeon Phi <u>co</u>-processor up to 61 Pentium cores
- Knights Landing & Knights Mill
  Xeon Phi full processor up to 36x dual-core Atom tiles

# Intel Knights Corner architecture

4





# The new Knights Landing architecture



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



### More details in a later set of slides...





### 公

### Variable Precision: What is VNNI-16?

- Vector Neural Network Instructions
- Variable precision
  - Inputs: 16-bit INT
  - Outputs: 32-bit INT
  - Semantics: 2 x int16 multiplies horizontally accumulating into single 32bit output
- Variable precision is best of both worlds
  - Same operations/instruction as 'half precision'
    - 2x OPS vs Single Precision
  - Similar output precision for optimal training convergence
    - 31 bits of INT32 vs 24 bits of mantissa in FP32
  - The obvious trade-off is the associated overhead on handling dynamic range in software (fixed precision)



AJProença, Advanced Architectures, MiEI, UMinho, 2017/18

(intel)



AJProença, Advanced Architectures, MiEI, UMinho, 2017/18

(intel)

### 公

### **Knights Mill Core**

### **Enhanced Knights Landing core**

- 2-way, OOO execution
- 4-way SMT
- 1 MB L2 bandwidth (64 bytes/cyc)
- 46 PA bits, 48 VA bits
- 2x 512b loads, 1x 512b store
- 32 KB D-cache (8-way), 32 KB I-cache (8-way)
- 72 inflight uops
- RS sizes: IEC (2 x 12), MEC (12), VPU (2 x 20)
- 1st level uTLB: 64 entries 2nd level dTLB: 256 4K, 128 2M, 16 1G pages



#### ISA: SSE, AVX, AVX512-F

DP stack

- 1 VPU port/core
- 1x 16 DP flops per cycle
- 6 cycles of latency

#### SP/VNNI stack

- 2 VPU ports/core
- 2 stacked FMAs per port
- 2x 64 SP flops per cycle
- 2x 128 VP ops/cycle
- 3+3 cycles of latency

AJProença, Advanced Architectures, MiEI, UMinho, 2017/18

(intel)

### PEZY-SC: <u>Peta</u> <u>Exa</u> <u>Z</u>etta <u>Y</u>otta-<u>S</u>uper<u>C</u>omputer: a 1024-core many-core processor chip



# **Beyond Vector/SIMD architectures**

### 公

### • 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
  - x86 many-core: Intel MIC / Xeon KNL
  - other many-core: ShenWei 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

• ...



### Sunway, or ShenWei, chip (<u>Chinese</u>:神威)

SW26010, 260 cores (in Sunway TaihuLight, #1 in TOP500 since June'16)

- 4x management cores (MPE) + 4x computer clusters (CPE)
- each CPE: 4x 64-core 64-bit RISC processors
- each core: only cache L1 & w/ 256-bit vector instructions
- next generation: SW52010, with 8x MPE-cores and 8x 64-cores CPE meshes



# **Beyond Vector/SIMD architectures**

### シ

### • 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

- x86 many-core: Intel MIC / Xeon KNL
- other many-core: ShenWei 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



# FPGA as a multiple configurable ISA



# FPGA as a computing accelerator



# **The Intel Programmable Acceleration Card**

### **PRODUCT BRIEF**



### Intel<sup>®</sup> Programmable Acceleration Card (Intel<sup>®</sup> PAC) with Intel<sup>®</sup> Arria<sup>®</sup> 10 GX FPGAs

#### Introduction

This PCIe-based FPGA acceleration card for data centers offers both inline and lookaside acceleration. It provides the performance and versatility of FPGA acceleration and is one of several platforms supported by the Acceleration Stack for Intel® Xeon® CPUs with FPGAs. This acceleration stack provides a common developer interface for both application and accelerator function developers, and includes drivers, application programming interfaces (APIs), and an FPGA interface manager. Together with acceleration libraries and development tools, the acceleration stack saves developer's time and enables code re-use across multiple Intel FPGA platforms. The card can be deployed in a variety of servers with its lowprofile form factor, low-power dissipation, and passive heat sink.

### **Targeted Workloads**

- Big data analytics
- Artificial intelligence
- Video transcoding
- Cyber security
- High-performance computing (HPC), such as genomics and oil and gas





### AJProença,

公入

• Financial technology, or FinTech

# Faster integration of

# programmable acceleration cards at Intel



# **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

- x86 many-core: Intel MIC / Xeon KNL
- other many-core: ShenWei 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

• ...

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

• ...

# **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"*



# *# cores/processing elements in several devices*



# Theoretical peak performance in several computing devices (DP)



AJProença, Advanced Architectures, MiEI, UMinho, 2017/18

25

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



AJProença, Advanced Architectures, MiEI, UMinho, 2017/18

# **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**



公入

29

Graphical Processing Units

SM

I-Cache MT Issue

C-Cache

SPUSP

SP

SP

SPUSP

SFU SFL

DP

Shared Memory

# **NVIDIA GPU Memory Structures**

- Each SIMD Lane has private section of off-chip DRAM
  - "Private memory" (Local Memory)
  - Contains stack frame, spilling registers, and private variables
- Each multithreaded SIMD processor (SM) also has local memory (Shared Memory)
  - Shared by SIMD lanes / threads within a block
- Memory shared by SIMD processors (SM) is GPU Memory, off-chip DRAM (Global Memory)
  - Host can read and write GPU memory



# 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

Warp Scheduler

Dispatch Unit

Instruction Cache

Warp Scheduler

Dispatch Unit

- Thus, two threads of SIMD instructions are scheduled every two clock cycles
- 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



# Fermi: Multithreading and Memory Hierarchy



公

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



公

# 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.

# Families in NVidia Tesla GPUs



# From Fermi into Kepler: The Memory Hierarchy



AJProença, Advanced Architectures, MiEI, UMinho, 2017/18

公





SMX 192 CUDA-cores

Ratio DPunit : SPunit -> 1 : 3

AJProença, Advanced Architectures, MiEI, UMinł

## From Fermi to Kepler core: SM and the SMX Architecture

								Ins	tructi	on Ca	che								
Warp Scheduler			Warp Scheduler				Warp Scheduler				Warp Scheduler								
Dispatch		h	Dispatch		Dispatch Dispatch			Dispatch Dispatch			Dispatch Dispatch								
							Regi	ster I	-ile ((	65.53	6 x 3	2-bit)							
Ŧ	Ŧ	Ŧ	+	Ŧ	÷	÷	_ <b>+</b> _	+	•		÷	+	+	÷	÷	÷	+	+	-
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	S
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	\$
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	s
Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	SFU	Core	Core	Core	DP Unit	Core	Core	Core	DP Unit	LD/ST	\$
-		200005						Inter	conne	ct Ne	twork								
							64 KB	Share B Bo	ed Me	emor	y 7 L1 ata C	Cac	ne						
	Tax		Tax			Tox	40 K	Tee	au-01				Ter			Tor		Tee	
	Tex		Tex			Tex		Tex			Tex		lex	( )		Tex		Tex	C



# The move from Kepler to Maxwell : from 15 SMXs to 48 SMMs in 6 GPCs



AJProença, Advanced Architectures, MiEI, UMinho, 2017/18

Dispatch

Unit Core Core Core DP Uni

.D/ST S

Register File (65,536 x 32-bit)

Core Core Core







#### Volta Architecture: 6x GPCs, 80 SMs SM Warp Scheduler (32 thread/clk Narp Scheduler (32 thread/clk Dispatch Unit (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 FP64 INT INT FP32 P32 FP3 FP64 INT INT FP64 INT INT FP64 INT INT P32 FP64 INT P32 FP64 INT INT FP64 INT P32 FP3 TENSOR TENSOR TENSOR TENSOR CORE CORE CORE CORE FP64 INT INT FP64 INT FP32 ED32 Volta SM: FP64 INT INT P32 FP3 FP64 INT FP32 64 CUDA-cores FP64 FP64 INT FP64 FP64 INT New: 8 Tensor-cores LD/ LD/ ST ST LD/ LD/ LD/ LD/ LD/ LD/ ST ST ST ST ST ST LD/ LD/ ST ST LD/ LD/ ST ST LD/ ST LD/ ST LD/ LD/ ST ST SFU SFU Ratio DPunit : SPunit -> 1 : 2 L0 Instruction Cache Warp Scheduler (32 throad/clk) Warp Scheduler (32 thread/clk) Dispatch Unit (32 thread/clk) Dispatch Unit (32 thread/clk) Register File (16,384 x 32-bit) Register File (16,384 x 32-bit) FP64 INT FP64 Volta V100 w/ 16GB HBM2 FP64 INT INT FP64 INT FP64 INT INT FP64 INT FP64 INT INT P32 FP32 FP64 INT P32 FP32 TENSOR TENSOR TENSOR TENSOR CORE CORE CORE CORE FP64 P32 P32 INT INT FP64 FP64 INT INT FP64 INT INT FP32 FP3 P32 FP3 DVIDIA FP64 FP64 INT INT P32 INT INT P32 FP3 FP64 INT INT FP64 INT INT 00 LD/ ST LD/ ST LD/ ST LD/ ST LD/ LD/ ST ST LD/ LD/ ST SFU SFU AJProenca, Advanced Architectures, MiEI, UM Tex Tex Tex Tex

l	Tesla V100	Tesla P100	Tesla M40	Tesla K40	Tesla Product	
Tesla a	GV100 (Volta)	GP100 (Pascal)	GM200 (Maxwell)	GK180 (Kepler)	GPU	
rece	80	56	24	15	SMs	
	40	28	24	15	TPCs	
	64	64	128	192	FP32 Cores / SM	
	5120	3584	3072	2880	FP32 Cores / GPU	
	32	32	4	64	FP64 Cores / SM	
	2560	1792	96	960	FP64 Cores / GPU	
	8	NA	NA	NA	Tensor Cores / SM	
	640	NA	NA	NA	Tensor Cores / GPU	
	1530 MHz	1480 MHz	1114 MHz	810/875 MHz	GPU Boost Clock	
	15.7	10.6	6.8	5.04	Peak FP32 TFLOP/s	
	7.8	5.3	.21	1.68	Peak FP64 TFLOP/s	
ANNOUN	125	NA	NA	NA	Peak Tensor Core TFLOP/s	
GIANT LEAP FUR	320	224	1 <b>92</b>	240	Texture Units	
21B xtors   TSMC 12n	4096-bit HBM2	4096-bit HBM2	384-bit GDDR5	384-bit GDDR5	Memory Interface	
7.5 FP64 TFLOPS   15	16 GB	16 GB	Up to 24 GB	Up to 12 GB	Memory Size	
NEW 120 Tensor TFLOF 20MB SM RE 1 16MB C	6144 KB	4096 KB	3072 KB	1536 KB	L2 Cache Size	
300 GB/s NVLink	Configurable up to 96 KB	64 KB	96 KB	16 KB/32 KB/48 KB	Shared Memory Size / SM	
	256KB	256 KB	256 KB	256 KB	Register File Size / SM	
	20480 KB	14336 KB	6144 KB	3840 KB	Register File Size / GPU	
	300 Watts	300 Watts	250 Watts	235 Watts	TDP	
https://devblogs.n	21.1 billion	15.3 billion	8 billion	7.1 billion	Transistors	
	815 mm²	610 mm <sup>2</sup>	601 mm <sup>2</sup>	551 mm²	GPU Die Size	
	12 nm FFN	16 nm FinFET+	28 nm	28 nm	Manufacturing Process	

# esla accelerators: recent evolution

### ANNOUNCING TESLA V100 GIANT LEAP FOR AI & HPC VOLTA WITH NEW TENSOR CORE

21B xtors | TSMC 12nm FFN | 815mm<sup>2</sup> 5,120 CUDA cores 2.5 FP64 TFLOPS | 15 FP32 TFLOPS IEW 120 Tensor TFLOPS 0MB SM RF | 16MB Cache | 16GBAB 00 GB/s NVLink

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

# **Beyond Vector/SIMD architectures**

### 公

### • 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
- other many-core: **ShenWei** 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
- focus on tensor/neural nets cores: NVidia, IBM, Intel NNP, Google TPU

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

• ...

# Machine learning w/ neural nets & deep learning...



### Key algorithms to train & classify use matrix products, but require lower precision numbers!

# NVidia Volta Architecture: the new Tensor Cores



Figure 8. Tensor Core 4x4 Matrix Multiply and Accumulate

For each SM: 8x 64 FMA ops/cycle 1k FLOPS/cycle!



Figure 9. Mixed Precision Multiply and Accumulate in Tensor Core

AJProença, Advanced Architectures, MiEI, UMinho, 2017/18

http://www.nvidia.com/content/gated-pdfs/Volta-Architecture-Whitepaper-v1.1.pdf

# NVidia competitors with neural net features: IBM TrueNorth chip array (August'2014)

rueNorth Chip 64 x 64 cores

### $\sim$

TrueNorth Chip:

- 4096 neurosynaptic cores Each core:
- 256 inputs (axons)
- 256 outputs (neurons)
- RAM w/ data for each neuron
- router (any neuron to any axon)





# NVidia competitors with neural net features: the IBM TrueNorth architecture



# NVidia competitors with neural net features: Intel Nervana Neural Network Processor, NNP

# History

ふ

- Nervana Engine announced in May'16
- Key features:
  - ASIC chip, focused on matrix multiplication, convolutions,... (for neural nets)
  - HBM2: 4x 8GB in-package storage & 1TB/sec memory access b/w
  - no h/w managed cache hierarchy (saves die area, higher compute density)
  - built-in networking (6 bi-directional high-b/w links)
  - separate pipelines for computation and data management
  - proprietary numeric format Flexpoint in-between floating point and fixed point precision
- Nervana acquired by Intel in August 2016:
  - renamed the project to "Lake Crest"
  - later to Nervana NNP, launched in October'17
  - Loihi test chip w/ self-learning capabilities announced in Sept'17, to be launched in 2018

AJProença, Advanced Architectures, MiEl, UMinho, 2017/18

oihi

# NVidia competitors with neural net features: Google Tensor Processing Unit, TPU (April'17)

### ふ

- The Matrix Unit: 65,536 (256x256) 8-bit multiply-accumulate units
- 700 MHz clock rate
- Peak: 92T operations/second
  - 65,536 \* 2 \* 700M
- >25X as many MACs vs GPU
- >100X as many MACs vs CPU
- 4 MiB of on-chip Accumulator memory
- 24 MiB of on-chip Unified Buffer, (activation memory)
- 3.5X as much on-chip memory vs GPU
- Two 2133MHz DDR3 DRAM channels
- 8 GiB of off-chip weight DRAM memory

# TPU: High-level Chip Architecture



# NVidia competitors with neural net features: Google Tensor Processing Unit, TPU (April'17)



RankBrain, StreetView & Google Translate

公

https://cloud.google.com/blog/big-data/2017/05/an-in-depth-look-at-googles-first-tensor-processing-unit-tpu

# NVidia competitors with neural net features: Google TPUv2 (September'17)



# **Beyond Vector/SIMD architectures**

### シ

### • 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

- x86 many-core: Intel MIC / Xeon KNL
- other many-core: **ShenWei** 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
- focus on tensor/neural nets cores: NVidia, IBM, Intel NNP, Google TPU

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

- <u>x86</u> multicore coupled with SIMT/SIMD cores: Intel i5/i7
- <u>ARMv8</u> cores coupled with SIMT/SIMD cores: NVidia Tegra

# Intel multicore coupled with GPU-cores



# **NVidia Tegra: SoC partnership with ARM** (1)



公

• Tegra 3 in Audi infotainment (2012) ...

The World's First Mobile Quad Core, Tegra 3 with 5<sup>th</sup> Companion Core for Low Power Quad Core, with 5th Companion Core CPU - Up to 1.4GHz Single Core, 1.3GHz Quad Core Up to 3x Higher GPU Performance GPU - 12 Core GeForce GPU Blu-Ray Quality Video VIDEO - 1080p High Profile @ 40Mbps Lower Power than Tegra 2 POWER Variable Symmetric Multiprocessing (vSMP) Up to 3x Higher Memory Bandwidth MEMORY — DDR3L-1500, LPDDR2-1066 IMAGING Up to 2x Faster ISP (Image Signal Processor) AUDIO HD Audio, 7.1 channel surround 2-6x Faster STORAGE - e.MMC 4.41, SD3.0, SATA-II

> Tegra 3 Nov'2011

AJProenca, Advanced Architectures, MiEI, UMinho, 2017/18

TERRER IN CORPORATION OF MALE AND AND

Tegra 4: replace the 32-bit ARM Cortex A9 by Cortex A15, and add 72 CUDA-cores



Tegra 4

May'2013 56

# **NVidia Tegra: SoC partnership with ARM** (2)

### 公

Replace the GPU block by 192 GPU-cores (*from Kepler*) and offer either 32/64-bit CPU cores => **Tegra K1** 



# **NVidia Tegra: SoC partnership with ARM (3)**

### Replace ... => Tegra K1



AJProença,

公入

# **NVidia Tegra: SoC partnership with ARM** (4)

### XX

 Replace the 5x 32-bit ARM by 2x4 32-bit Cortex (A57 & A53) and the 192 Kepler CUDA cores by 256 Maxwell => Tegra X1



# NVidia Tegra: pathway towards ARM-64 (1)

### 公

 Upgrade 32-bit ARM to 64-bit ARMv8 (*Denver 2 & A57*) and replace Maxwell cores by Pascal ones => Parker Aug'2016



### "PARKER" CPU COMPLEX

- 2x Denver2 + 4x Cortex-A57
- Fully Coherent HMP system
  - Proprietary Coherent Interconnect
- ARM V8 64-bit
- Highest performance ARM CPU
  - 2nd generation Denver core
  - Significant Perf/W improvements
- Dynamic Code Optimization
  - OoO execution without the power
  - Optimize once, use many times
- 7-wide superscalar
- Low power retention states



# NVidia Tegra: pathway towards ARM-64 (2)

### 公

 Increment #ARMv8-cores (*custom architecture*) and replace Pascal cores by Volta (w/ tensor cores) => Xavier Jan'2018?

E CVA BK VP		NVID	IA ARM SoCs	
8 LO CORE		Xavier	Parker	Erista (Tegra XI)
512 CORE GPU CPU	CPU	8x NVIDIA Custom ARM	2x NVIDIA Denver + 4x ARM Cortex-A57	4x ARM Cortex-A57 + 4x ARM Cortex-A53
	GPU	Volta, 512 CUDA Cores	Pascal, 256 CUDA Cores	Maxwell, 256 CUDA Cores
	Memory	?	LPDDR4, 128-bit Bus	LPDDR3, 64-bit Bus
	Video Processing	7680x4320 Encode & Decode	3840x2160p60 Decode 3840x2160p60 Encode	3840x2160p60 Decode 3840x2160p30 Encode
	Transistors	7B	?	?
	Manufacturing Process	TSMC 16nm FinFET+	TSMC 16nm FinFET+	TSMC 20nm Planar

# **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
    - other many-core: ShenWei 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
    - focus on tensor/neural nets cores: NVidia, IBM, Intel NNP, Google TPU

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

- <u>x86</u> multicore coupled with SIMT/SIMD cores: **Intel** i5/i7
- <u>ARMv8</u> cores coupled with SIMT/SIMD cores: **NVidia** Tegra