# Master Informatics Eng.

2016/17

*A.J.Proença*

## Data Parallelism 2 (*SIMD++, Intel MIC*)
### *(most slides are borrowed)*
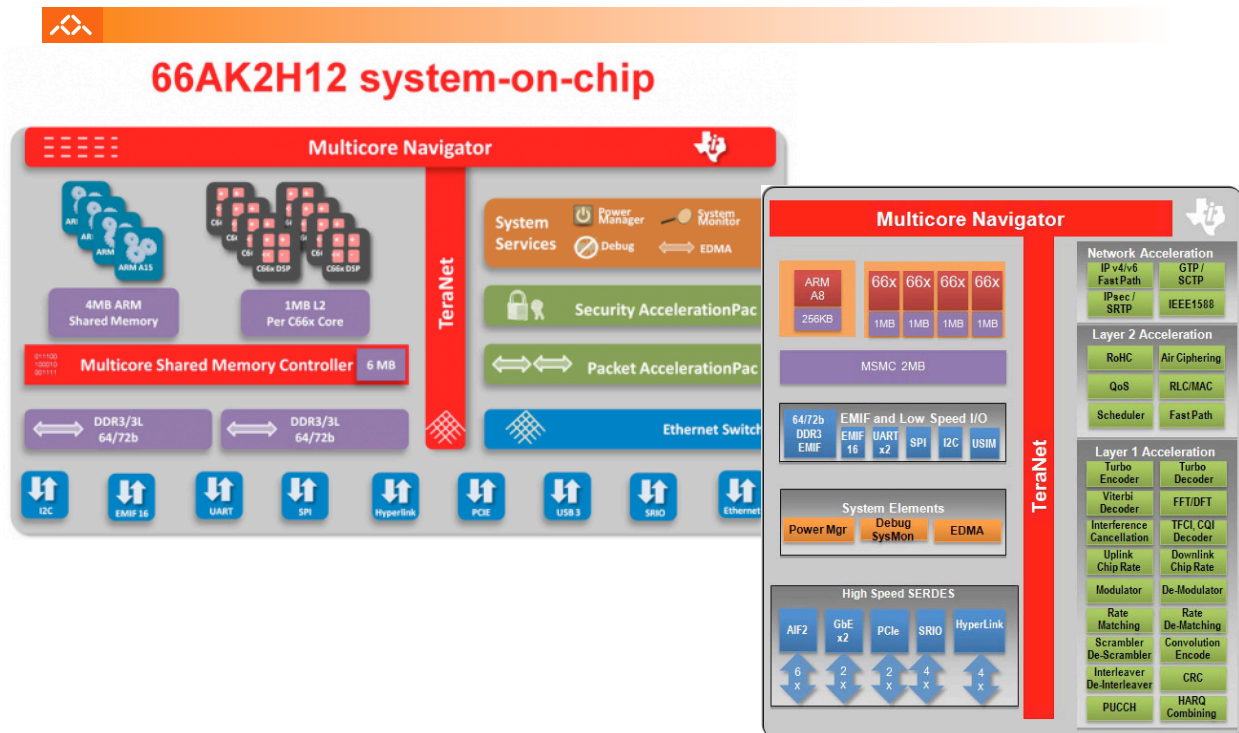
---

# *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

  - **CPU cores with wider vectors and/or SIMD cores**:
    - DSP VLIW cores with vector capabilities: **Texas Instruments** (...?)
    - PPC cores coupled with SIMD cores: **Cell** (past...) , **IBM Power BQC...**
    - ARM64 cores coupled with SIMD cores: from Tegra to Parker (**NVidia**) (...?)
    - x86 many-core: **Intel** MIC / Xeon KNL, **AMD** FirePro**...**
    - other many-core: **ShenWay** 260, **Adapteva** Epiphany-V**...**

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

# Texas Instruments: Keystone DSP architecture

# 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
  - **CPU cores with wider vectors and/or SIMD cores**:
    - DSP VLIW cores with vector capabilities: **Texas Instruments** (...?)
    - PPC cores coupled with SIMD cores: **Cell** (past...) , **IBM Power BQC...**
    - ARM64 cores coupled with SIMD cores: from Tegra to Parker (**NVidia**) (...?)
    - x86 many-core: **Intel** MIC / Xeon KNL, **AMD** FirePro**...**
    - other many-core: **ShenWay** 260, **Adapteva** Epiphany-V**...**
  - **coprocessors (require a host scalar processor): accelerator devices**
    - on disjoint physical memories (e.g., **Xeon KNC** with PCI-Expr, **PEZY-SC**)
    - focus on SIMT/SIMD to hide memory latency: **GPU**-type approach
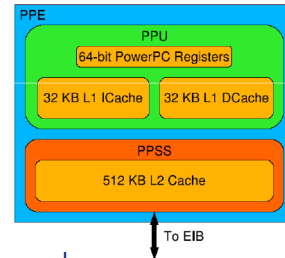    - ISA-free architectures, code compiled to silica: **FPGA**

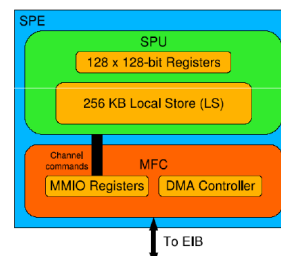# IBM Cell Broadband Engine (PPE)

- Heterogeneous multicore processor

  - 1 x Power Processor Element (PPE)

    - 64-bit Power-architecture-compliant processor

    - Dual-issue, in-order execution, 2-way SMT processor

    - PowerPC Processor Unit (PPU)

      - 32 KB L1 IC, 32 KB L1 DC, VMX unit

    - PowerPC Processor Storage Subsystem (PPSS)

      - 512 KB L2 Cache

    - General-purpose processor to run OS and control-intensive code

    - Coordinates the tasks performed by the remaining cores
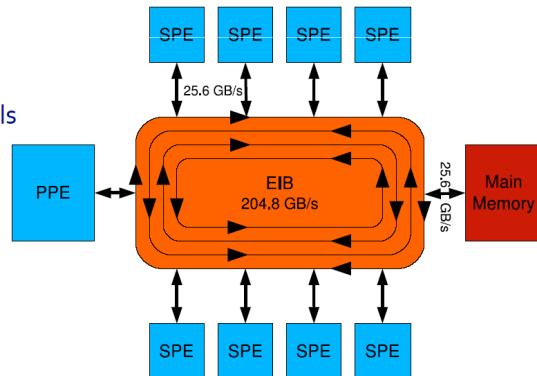
# IBM Cell Broadband Engine (SPE)

- Heterogeneous multicore processor

  - 8 x Synergistic Processing Element (SPE)

    - Dual-issue, in-order execution, 128-bit SIMD processors

    - Synergistic Processor Unit (SPU)

      - SIMD ISA (four different granularities)

      - 128 x 128-bit SIMD register file

      - **256 KB Local Storage (LS) for code/data**

    - Memory Flow Controller (MFC)

      - Memory-mapped I/O registers (MMIO Registers)

      - DMA Controller: commands to transfer data in and out

    - Custom processors specifically designed for data-intensive code

    - Provide the main computing power of the Cell BE

# IBM Cell Broadband Engine (*EIB*)

- Element Interconnect Bus (EIB)
  - Interconnects PPE, SPEs, and the memory and I/O interface controllers
    - 4 x 16 Byte-wide rings (2 clockwise and 2 counterclockwise)
  - Up to three simultaneous data transfers per ring
  - Shortest path algorithm for transfers
- Memory Interface Controller (MIC)
  - 2 x Rambus XDR I/O memory channels

  (accesses on each channel
  of 1-8, 16, 32, 64 or 128 Bytes)
- Cell BE Interface (BEI)
  - 2 x Rambus FlexIO I/O channels

# IBM Cell Broadband Engine (chip)

**Architecture**

# IBM Power BlueGene/Q Compute *(chip)*
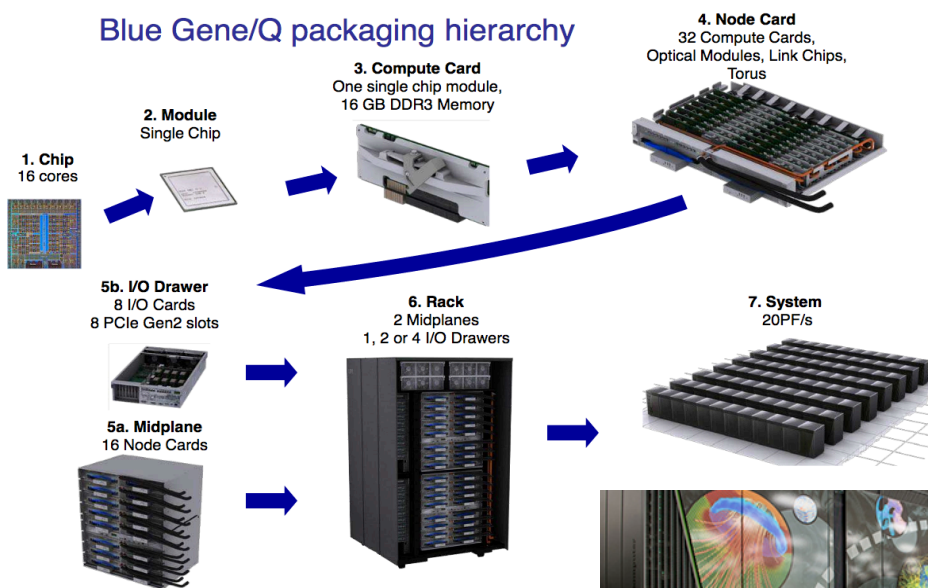


## Features:

- launched in 2010/11
  (TOP500: #1 in Jun12, #4 in Jun16)

- 18-cores   (16 compute,
       1 OS support, 1 redundant)
  - each 4-way  multi-threaded
  - 64 bits PowerISA
  - 1.6   GHz
  - L1   I/D cache = 16  kB/16  kB
  - each core has Quad FPU
    (4-wide double precision SIMD)

- shared L2 cache: 32 MB

- dual memory controller

*AJProença, Advanced Architectures, MiEI, UMinho, 2016/17*                                           *9*

# IBM Power BlueGene/Q Compute *(Sequoia system)*

## Blue Gene/Q packaging hierarchy



**1. Chip**
16 cores

**2. Module**
Single Chip

**3. Compute Card**
One single chip module,
16 GB DDR3 Memory

**4. Node Card**
32 Compute Cards,
Optical Modules, Link Chips,
Torus

**5b. I/O Drawer**
8 I/O Cards
8 PCIe Gen2 slots

**5a. Midplane**
16 Node Cards

**6. Rack**
2 Midplanes
1, 2 or 4 I/O Drawers

**7. System**
20PF/s

Ref: SC2010

*AJProença, Advanced Architectures, MiEI, UMinho, 2016/17*

# 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
  - **CPU cores with wider vectors and/or SIMD cores**:
    - DSP VLIW cores with vector capabilities: **Texas Instruments** (...?)
    - PPC cores coupled with SIMD cores: **Cell** (past...) , **IBM Power BQC...**
    - ARM64 cores coupled with SIMD cores: from Tegra to Parker (**NVidia**) (...?)
    - x86 many-core: **Intel** MIC / Xeon KNL**, AMD** FirePro**...**
    - other many-core: **ShenWay** 260, **Adapteva** Epiphany-V**...**
  - **coprocessors (require a host scalar processor): accelerator devices**
    - on disjoint physical memories (e.g., **Xeon KNC** with PCI-Expr, **PEZY-SC**)
    - focus on SIMT/SIMD to hide memory latency: **GPU**-type approach
    - ISA-free architectures, code compiled to silica: **FPGA**

*AJProença, Advanced Architectures, MiEI, UMinho, 2016/17*                  *11*


# NVidia: pathway towards ARM-64 (1)
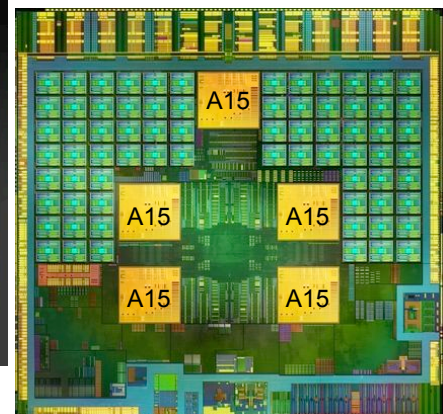
- Pick a successful line:
  Tegra 3, 4, ...

- Replace the 32-bit ARM Cortex A9 by Cortex A15, and add 72 CUDA-
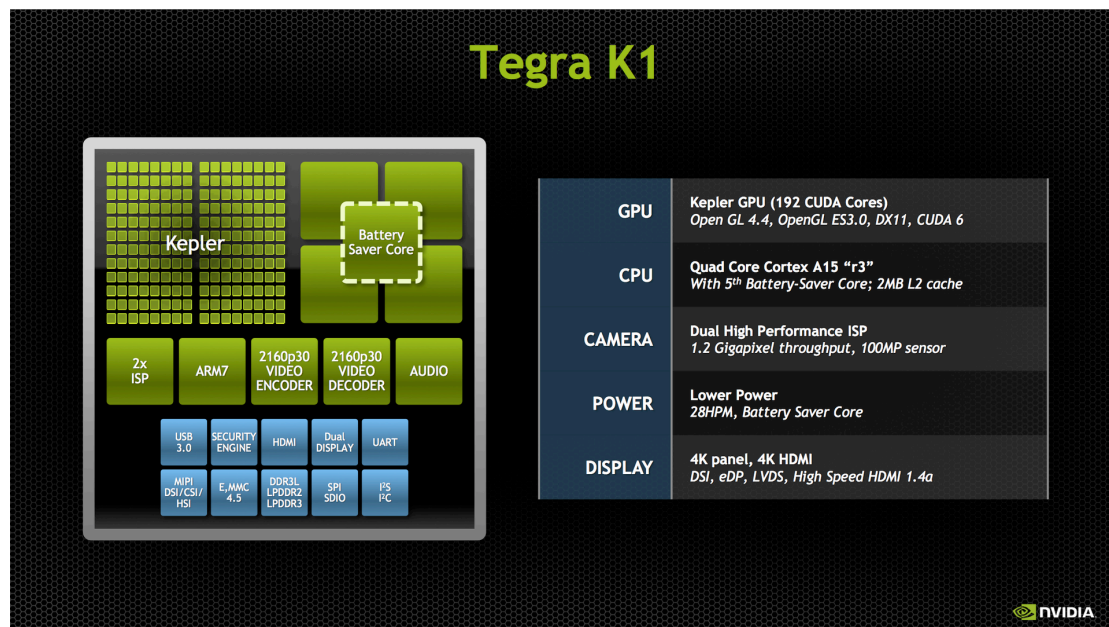


**Tegra 3**



**Tegra 4**

*AJProença, Advanced Architectures, MiEI, UMinho, 2016/17*                  *12*

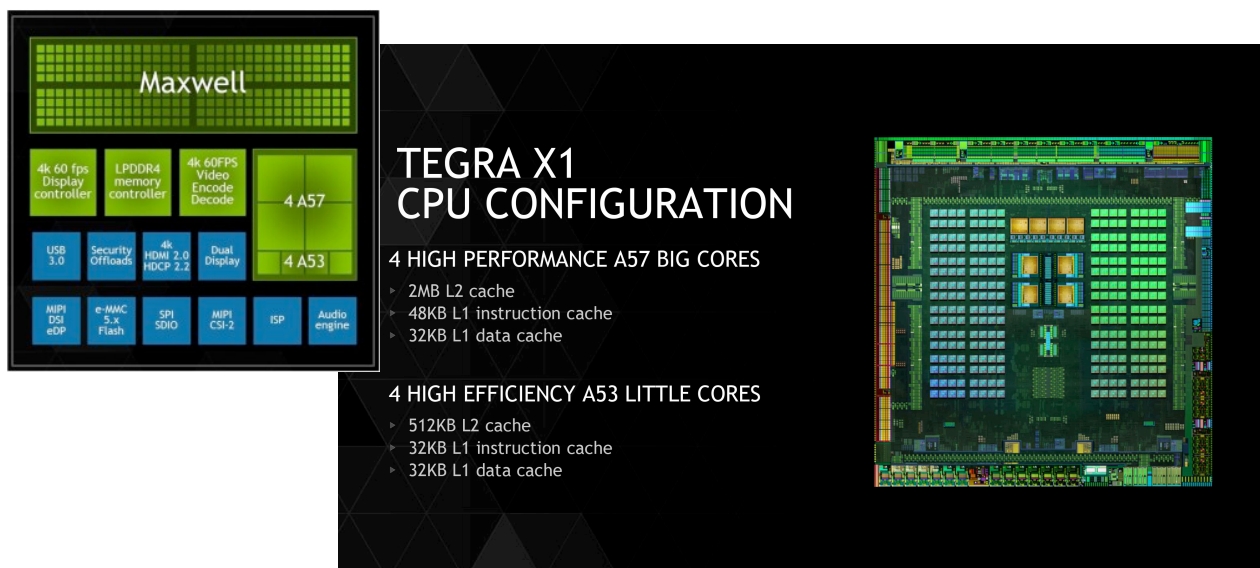## NVidia: pathway towards ARM-64 *(2)*

- Replace the GPU block by 192 GPU-cores (*from Kepler*) and keep the 5x 32-bit CPU cores (*Cortex A15*) => **Tegra K1**



## NVidia: pathway towards ARM-64 *(3)*

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

- Upgrade 32-bit ARM to 32- & 64-bit ARM (*Denver 2*) and replace the Maxwell CUDA cores by Pascal ones => **Parker**

## TEGRA KEY FEATURE EVOLUTION
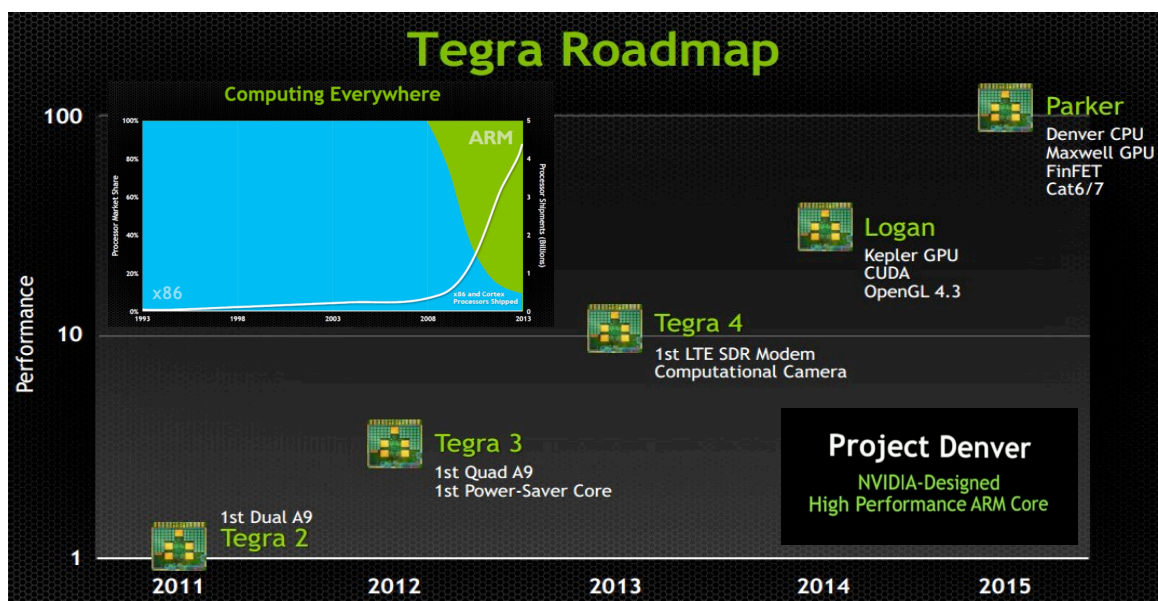
|  | TK1 | TX1 | "PARKER" |
|---|---|---|---|
| GPU | Kepler, 192 CUDA cores | Maxwell, 256 CUDA cores | Pascal, 256 CUDA cores |
| CPU | 4+1 A15, 2MB+512K L2 ARM v7 32b Or 2 Denver 1, 2MB L2 64b | 4x A57 2MB L2 + 4x A53 512KB L2 ARM v8 64b | 2x Denver 2 2MB L2 + 4x A57 2MB L2 ARM v8 64b Coherent HMP Architecture |
| Camera | 4 cameras | 6 cameras | Auto HDR 12 cameras |
| Memory | 64b LPDDR2/3, DDR3L 15 GB/s (LP3, DDR3L) | 64b LPDDR4, 25GB/s | 128b LPDDR4, 50 GB/s, ECC |
| Display | Dual Pipeline 4K@30fps 24bpp | Dual Pipeline 4K@60fps | Triple Pipeline 4K@60fps |

Tegra Roadmap

# 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
  - **CPU cores with wider vectors and/or SIMD cores**:
    - DSP VLIW cores with vector capabilities: **Texas Instruments** (...?)
    - PPC cores coupled with SIMD cores: **Cell** (past...) , **IBM Power BQC...**
    - ARM64 cores coupled with SIMD cores: from Tegra to Parker (**NVidia**) (...?)
    - **x86** many-core: **Intel MIC / Xeon KNL, AMD FirePro...**
    - other many-core: **ShenWay** 260, **Adapteva** Epiphany-V**...**
  - **coprocessors (require a host scalar processor): accelerator devices**
    - on disjoint physical memories (e.g., **Xeon KNC** with PCI-Expr, **PEZY-SC**)
    - focus on SIMT/SIMD to hide memory latency: **GPU**-type approach
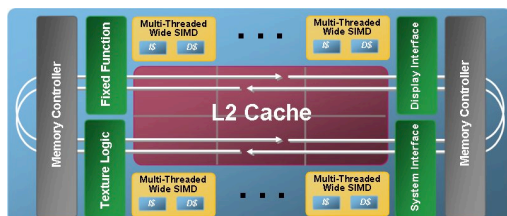    - ISA-free architectures, code compiled to silica: **FPGA**

*AJProença, Advanced Architectures, MiEI, UMinho, 2016/17*      *17*
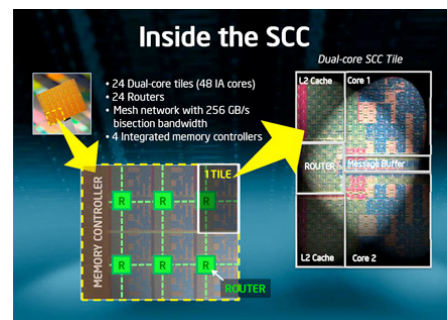
---

# Intel MIC: Many Integrated Core

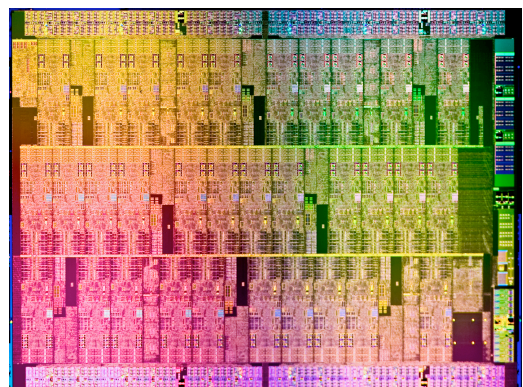## Intel evolution, from:
- Larrabee (80-core GPU)    &   SCC



**S**ingle-chip
**C**loud
**C**omputer,
24x
dual-core tiles

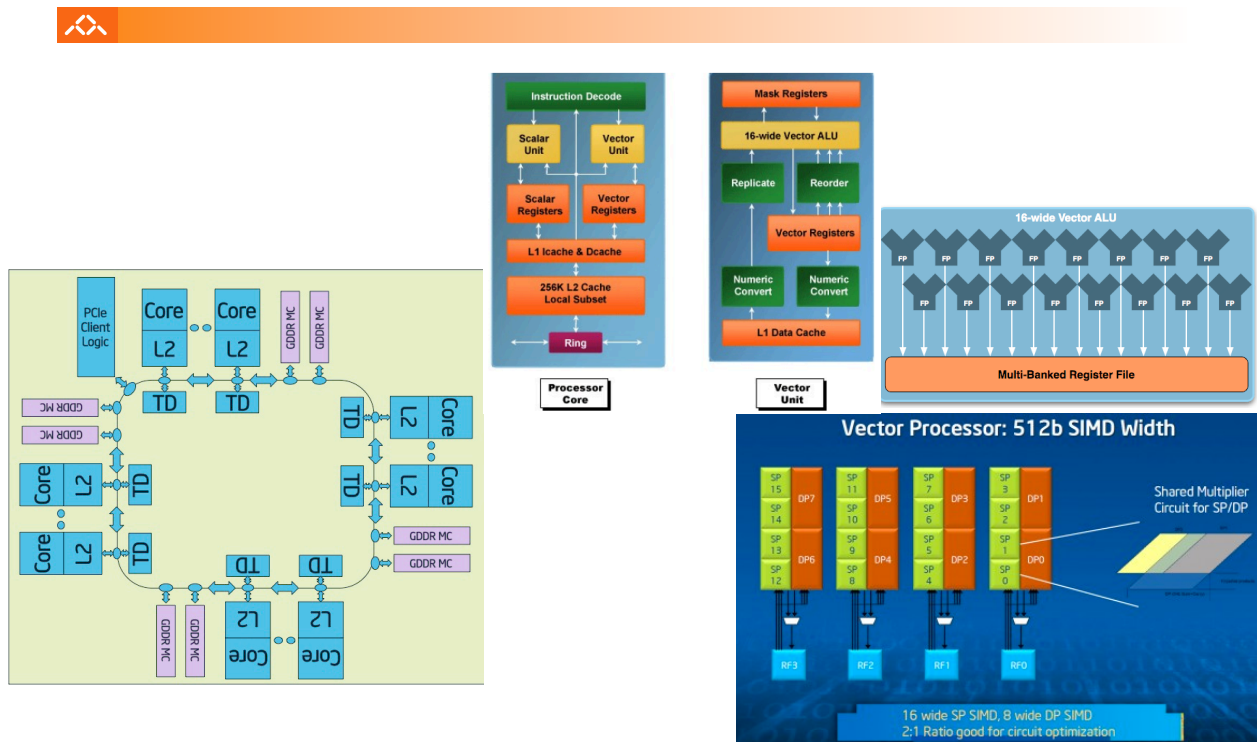## to MIC:
- Knights Ferry (pre-production, Stampede)
- Knights Corner
  Xeon Phi co-processor up to 61 Pentium cores
- Knights Landing
  Xeon Phi full processor,
  36x dual-core tiles with 64-bit Atoms
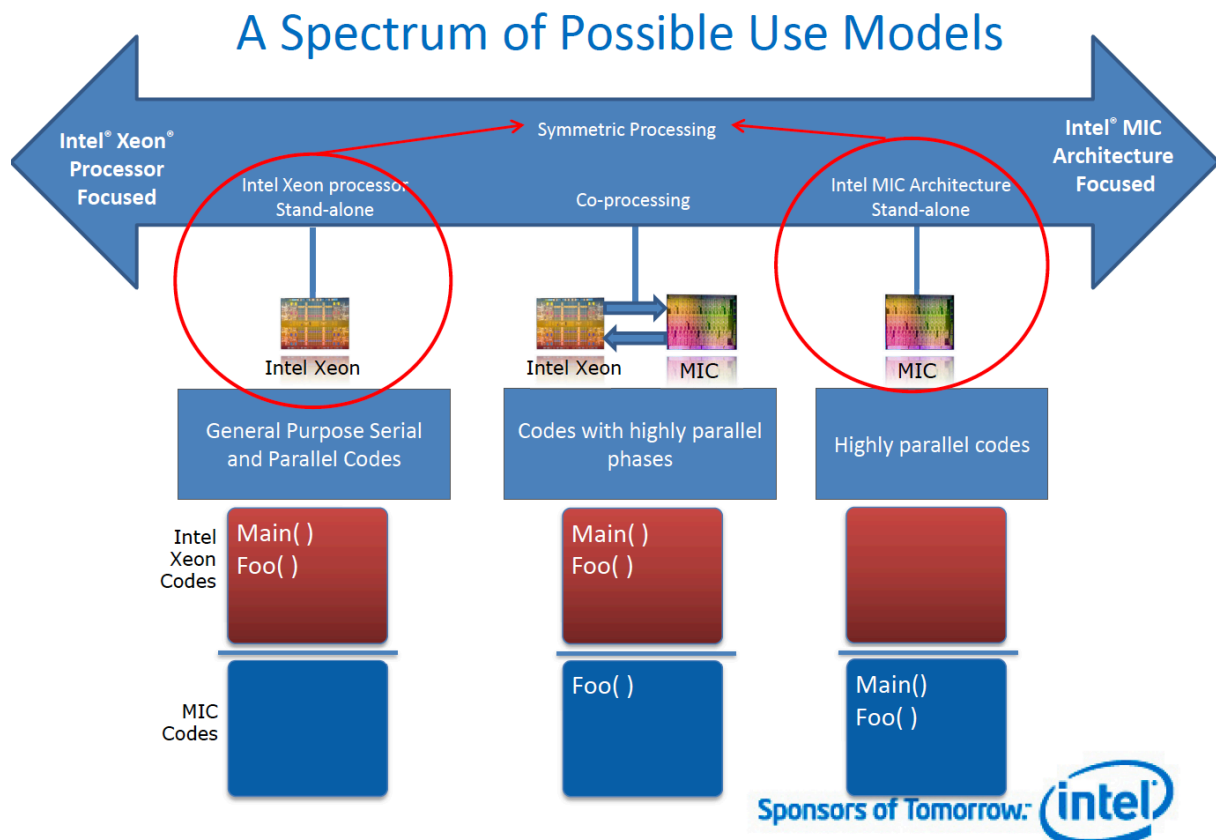
*AJProença, Advanced Architectures, MiEI, UMinho, 2016/17*

# Intel Knights Corner architecture

# A Spectrum of Possible Use Models

# The new Knights Landing architecture



## Innovation

### High-bandwidth In-Package Memory

**Performance for memory-bound workloads**

**Flexible memory usage models**

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



### Knights Landing Overview

**TILE:** 2 VPU | CHA / 1MB L2 | 2 VPU ; Core | | Core

**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)

**Vector Peak Perf:** 3+TF DP and 6+TF SP Flops
**Scalar Perf:** ~3x over Knights Corner
**Streams Triad (GB/s):** MCDRAM : 400+; DDR: 90+

36 Tiles connected by 2D Mesh Interconnect

Omni-path not shown

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

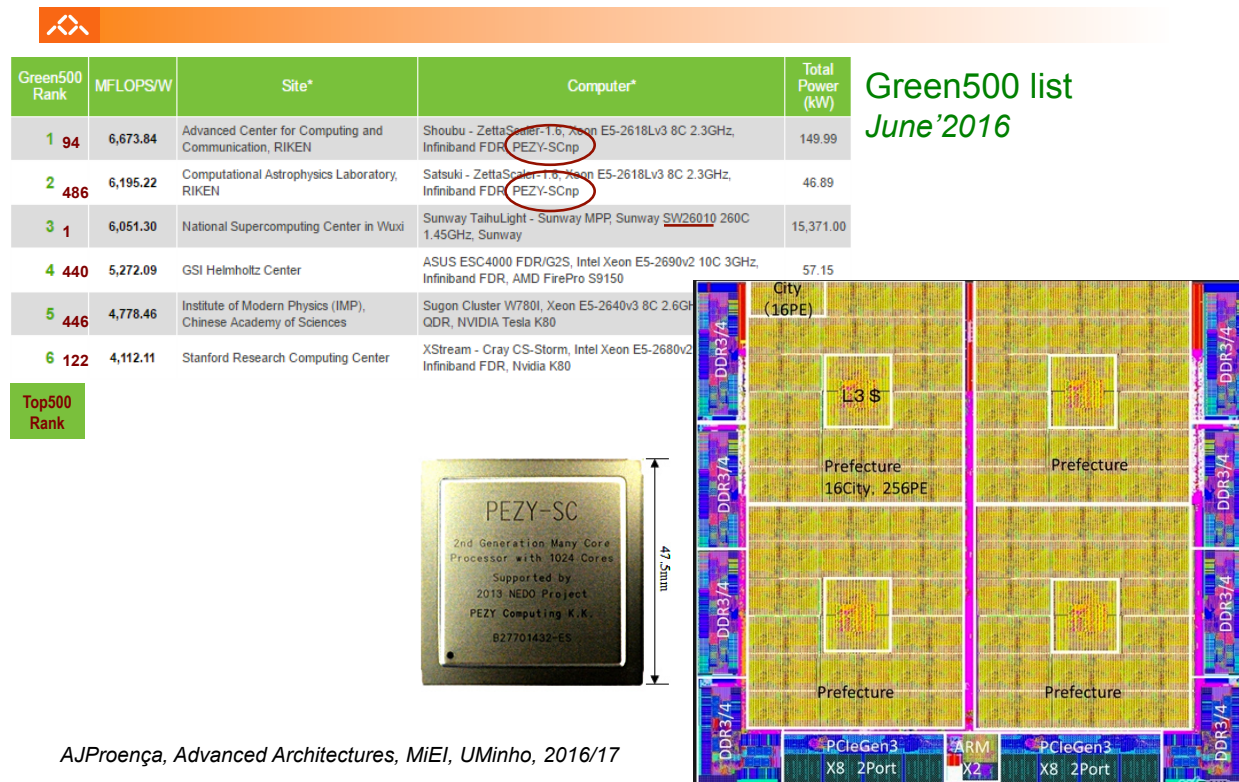| Green500 Rank | MFLOPS/W | Site* | Computer* | Total Power (kW) |
|---|---|---|---|---|
| 1  94 | 6,673.84 | Advanced Center for Computing and Communication, RIKEN | Shoubu - ZettaScaler-1.6, Xeon E5-2618Lv3 8C 2.3GHz, Infiniband FDR, PEZY-SCnp | 149.99 |
| 2  486 | 6,195.22 | Computational Astrophysics Laboratory, RIKEN | Satsuki - ZettaScaler-1.6, Xeon 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 |
| 5  446 | 4,778.46 | Institute of Modern Physics (IMP), Chinese Academy of Sciences | Sugon Cluster W780I, Xeon E5-2640v3 8C 2.6GHz, QDR, NVIDIA Tesla K80 | |
| 6  122 | 4,112.11 | Stanford Research Computing Center | XStream - Cray CS-Storm, Intel Xeon E5-2680v2, Infiniband FDR, Nvidia K80 | |

**Top500 Rank**

Green500 list
*June'2016*

*AJProença, Advanced Architectures, MiEI, UMinho, 2016/17*



PEZY-SC
2nd Generation Many Core Processor with 1024 Cores
Supported by 2013 NEDO Project
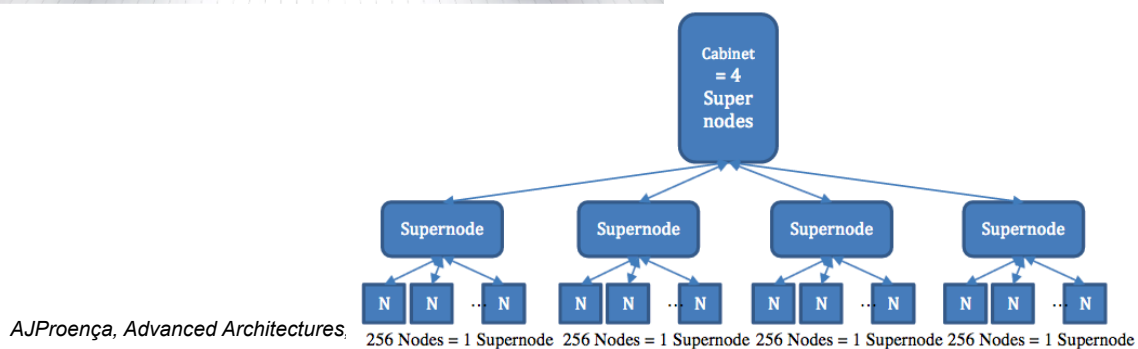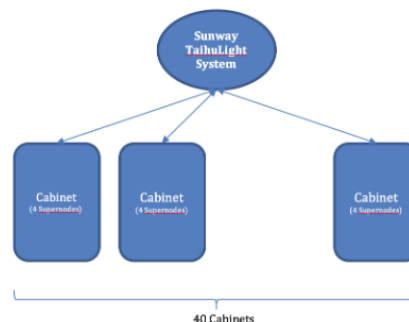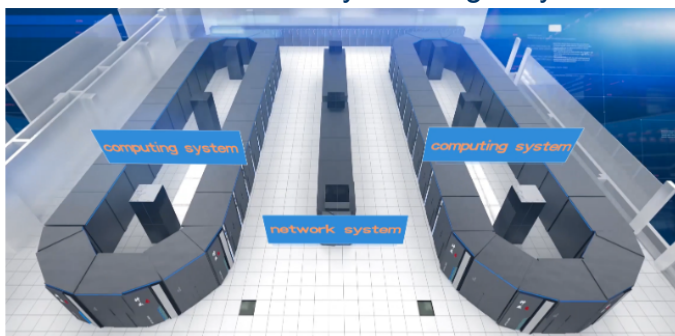PEZY Computing K.K.
B27701432-ES

47.5mm

---

# 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

  - **CPU cores with wider vectors and/or SIMD cores**:
    - DSP VLIW cores with vector capabilities: **Texas Instruments** (...?)
    - PPC cores coupled with SIMD cores: **Cell** (past...) , **IBM Power BQC...**
    - ARM64 cores coupled with SIMD cores: from Tegra to Parker (**NVidia**) (...?)
    - x86 many-core: **Intel** MIC / Xeon KNL, **AMD** FirePro...
    - other many-core: **ShenWay** 260, **Adapteva** Epiphany-V...

  - **coprocessors (require a host scalar processor): accelerator devices**
    - on disjoint physical memories (e.g., **Xeon KNC** with PCI-Expr, **PEZY-SC**)
    - focus on SIMT/SIMD to hide memory latency: **GPU**-type approach
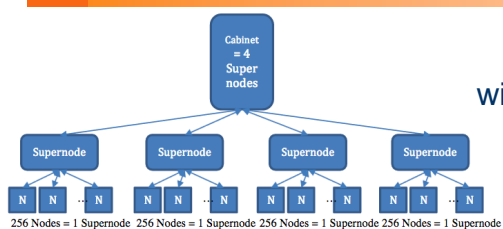    - ISA-free architectures, code compiled to silica: **FPGA**
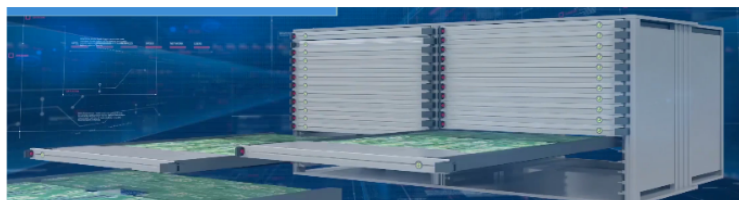
*#1 in June'16 TOP500:*
*Sunway TaihuLight*

## Overview of the Sunway TaihuLight System



256 Nodes = 1 Supernode  256 Nodes = 1 Supernode  256 Nodes = 1 Supernode  256 Nodes = 1 Supernode

*AJProença, Advanced Architectures,*

---



*#1 in June'16 TOP500:*
*Sunway TaihuLight*

256 Nodes = 1 Supernode   256 Nodes = 1 Supernode   256 Nodes = 1 Supernode   256 Nodes = 1 Supernode

One cabinet
with 4 Supernodes

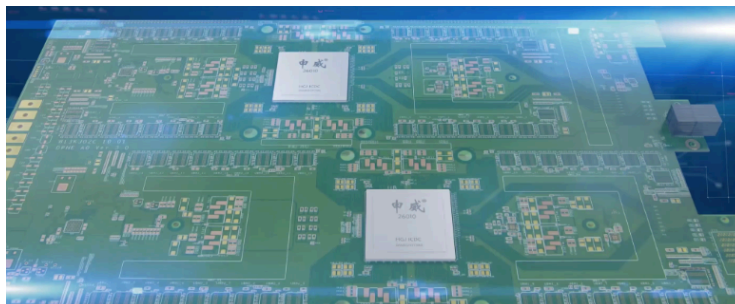One Supernode
with 32 boards

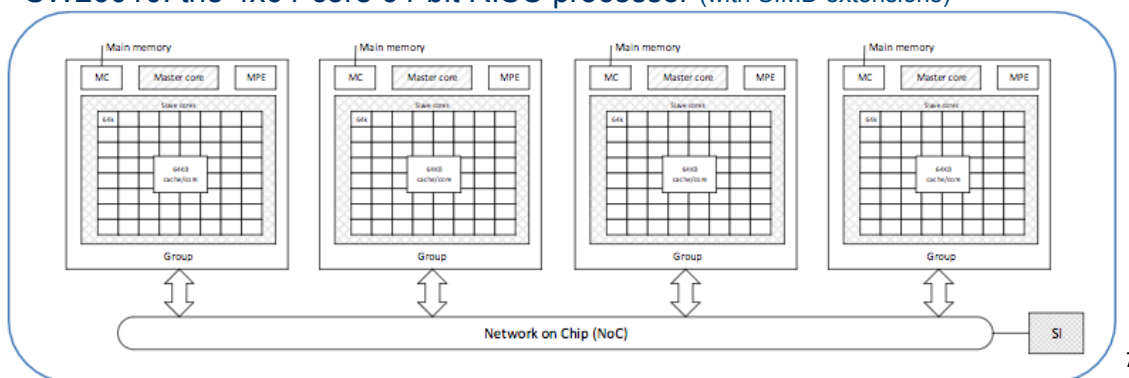One board with 4 cards,
2 up & 2 down

*#1 in June'16 TOP500:*
*Sunway TaihuLight*

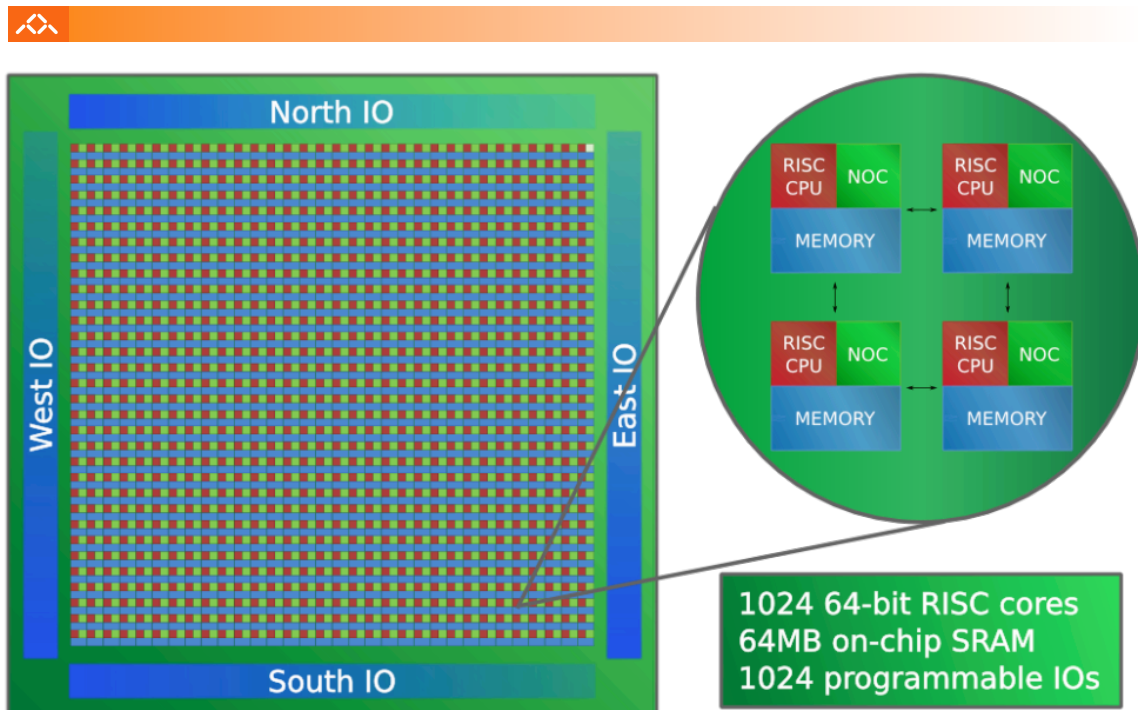One card with two nodes
*(two SW26010 chips)*



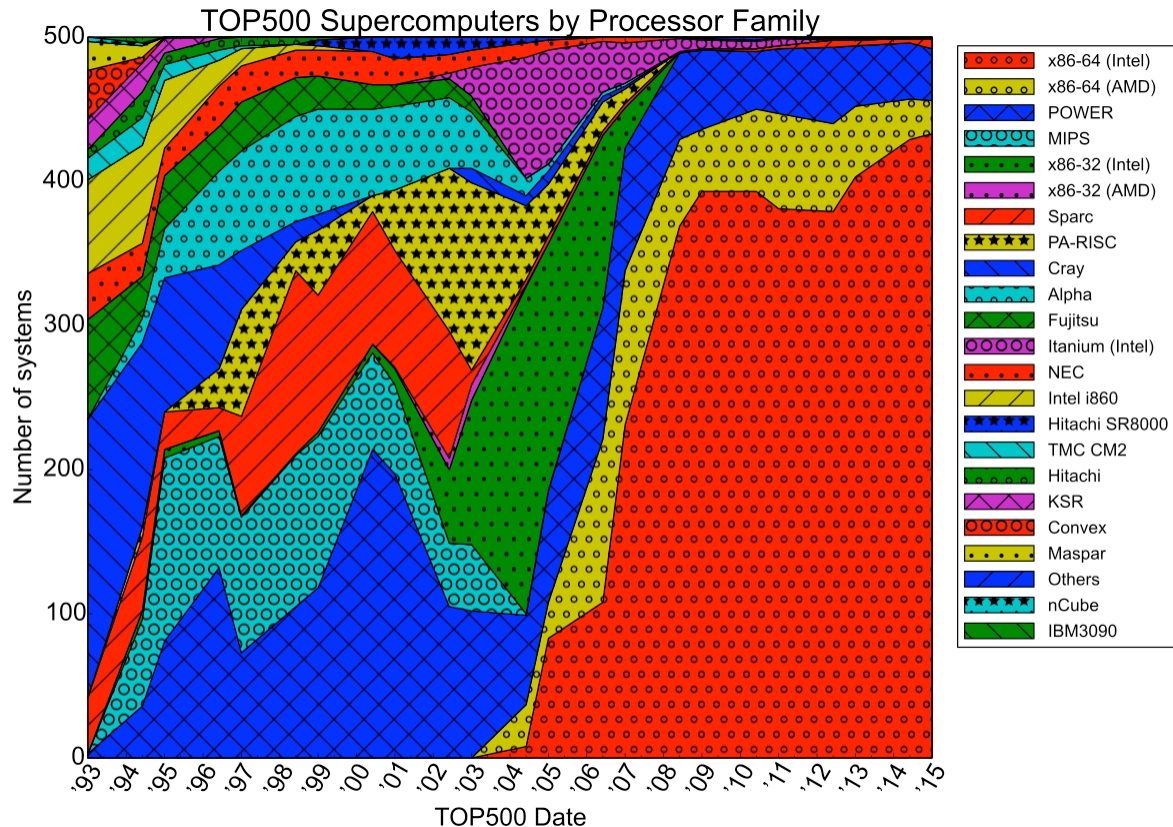SW26010: the 4x64-core 64-bit RISC processor (with SIMD extensions)

7

*Adapteva announcement in Oct'16:*
*Epiphany-V, a 1024-core RISC chip*



1024 64-bit RISC cores
64MB on-chip SRAM
1024 programmable IOs

# *Top500: Processor family distribution over all systems*



TOP500 Supercomputers by Processor Family

Legend:
- x86-64 (Intel)
- x86-64 (AMD)
- POWER
- MIPS
- x86-32 (Intel)
- x86-32 (AMD)
- Sparc
- PA-RISC
- Cray
- Alpha
- Fujitsu
- Itanium (Intel)
- NEC
- Intel i860
- Hitachi SR8000
- TMC CM2
- Hitachi
- KSR
- Convex
- Maspar
- Others
- nCube
- IBM3090

Axes: Number of systems (0–500) vs TOP500 Date ('93–'15)

---

# *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
  - **CPU cores with wider vectors and/or SIMD cores**:
    - DSP VLIW cores with vector capabilities: **Texas Instruments** (...?)
    - PPC cores coupled with SIMD cores: **Cell** (past...) , **IBM Power BQC...**
    - ARM64 cores coupled with SIMD cores: from Tegra to Parker (**NVidia**) (...?)
    - x86 many-core: **Intel** MIC / Xeon KNL**, AMD** FirePro**...**
    - other many-core: **ShenWay** 260, **Adapteva** Epiphany-V**...**
  - **coprocessors (require a host scalar processor): accelerator devices**
    - on disjoint physical memories (e.g., Xeon KNC with **PCI-E**xpr, **PEZY-SC**)
    - focus on SIMT/SIMD to hide memory latency: **GPU**-type approach
    - ISA-free architectures, code compiled to silica: **FPGA**