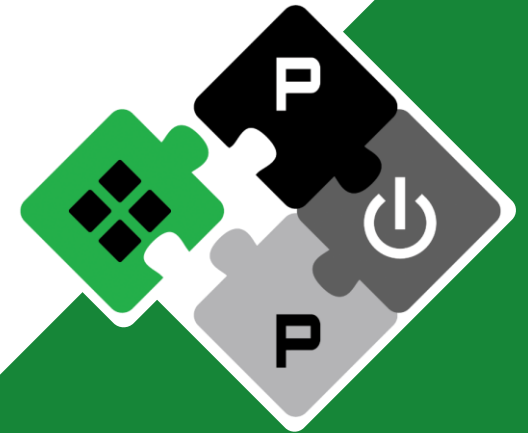# Enabling Sustainable AI

## an Open Computing Platform Perspective

Luca Benini

lbenini@ethz.ch, luca.Benini@unibo.it

**PULP Platform**
Open Source Hardware, the way it should be!
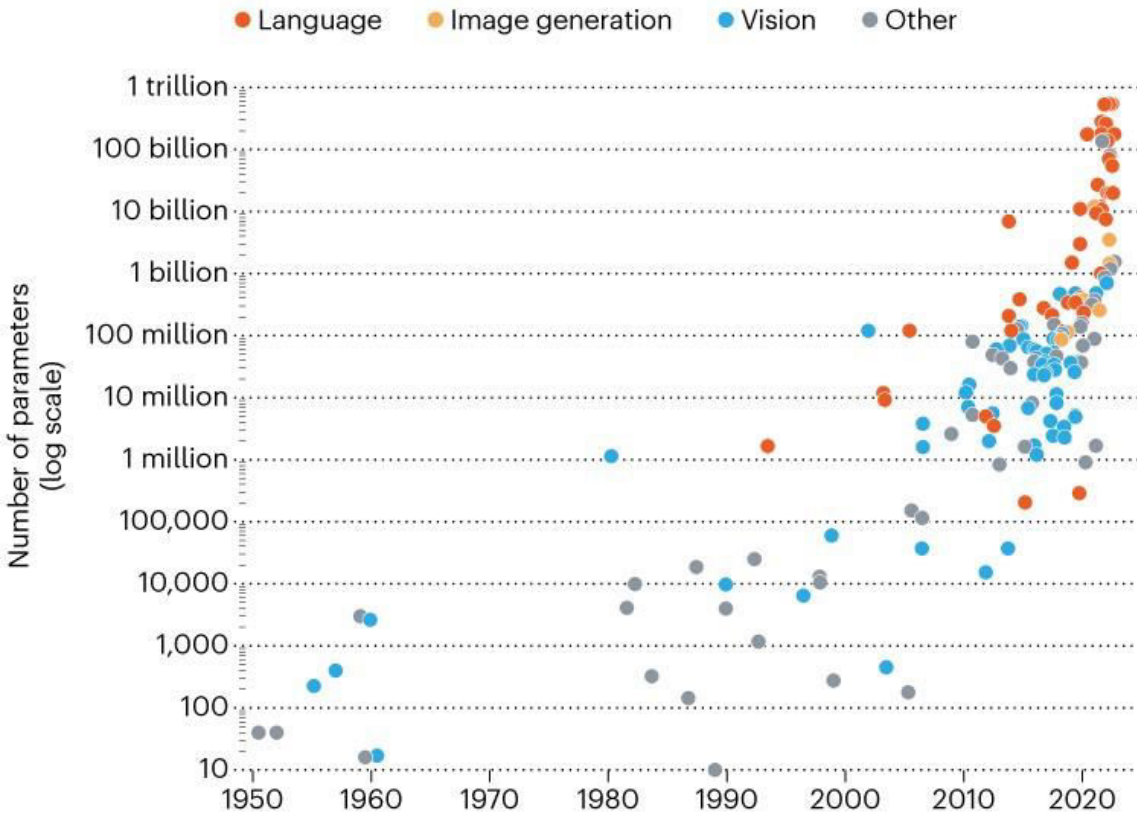
@pulp_platform

pulp-platform.org

youtube.com/pulp_platform

# AI is Power Bound from Cloud to Edge

**Datacenter PMAX < 150MW**
**On-car Computing PMAX < 1.5KW**

## THE DRIVE TO BIGGER AI MODELS

The scale of artificial-intelligence neural networks is growing exponentially, as measured by the models' parameters (roughly, the number of connections between their neurons)*.
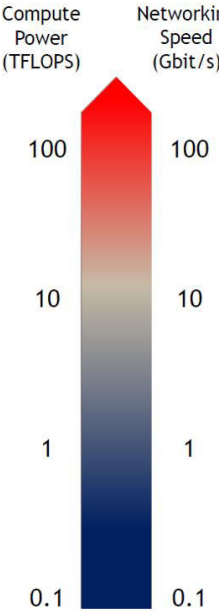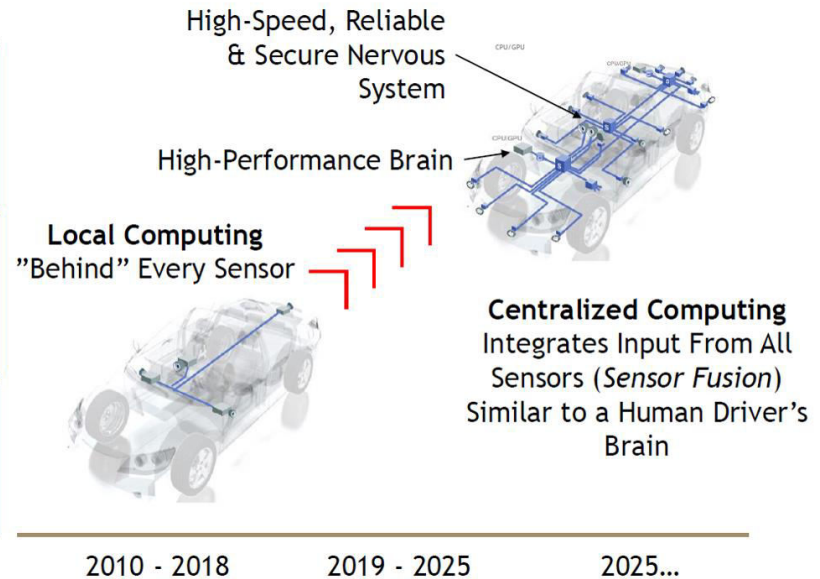


- Language
- Image generation
- Vision
- Other

Number of parameters (log scale): 1 trillion, 100 billion, 10 billion, 1 billion, 100 million, 10 million, 1 million, 100,000, 10,000, 1,000, 100, 10

Years: 1950, 1960, 1970, 1980, 1990, 2000, 2010, 2020

[Nature'23]

**10x every 2 years**

## Path Towards Full Autonomy



High-Speed, Reliable & Secure Nervous System

High-Performance Brain

**Local Computing** "Behind" Every Sensor

**Centralized Computing** Integrates Input From All Sensors (*Sensor Fusion*) Similar to a Human Driver's Brain

Level 4-5 Self Driving
Level 2-3 Decision Assistant
Level 1-2 Simple Aid

[SCR'23]

2010 - 2018    2019 - 2025    2025...

Compute Power (TFLOPS): 100, 10, 1, 0.1
Networking Speed (Gbit/s): 100, 10, 1, 0.1

**Energy Efficiency from Scaling**

$$\left(\frac{1}{Power \cdot Time}\right)$$

**10x every 12 years...**

ETH zürich    ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA

# Energy-Efficient Computing: Core to Platform
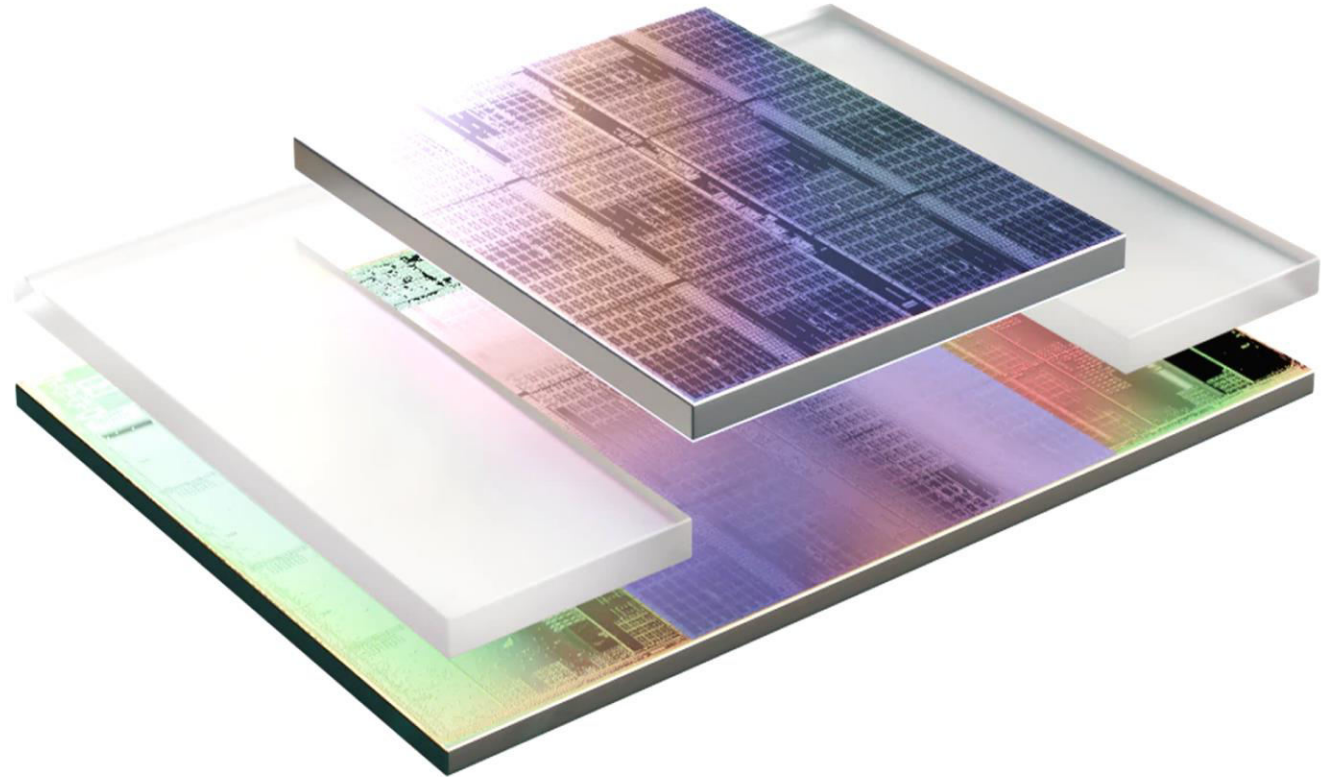
- **Processing Element (PE)**
  - Efficient operations and local storage

- **Chip**
  - Data movement and storage hierarchy

- **Full platform**
  - IOs, main memory
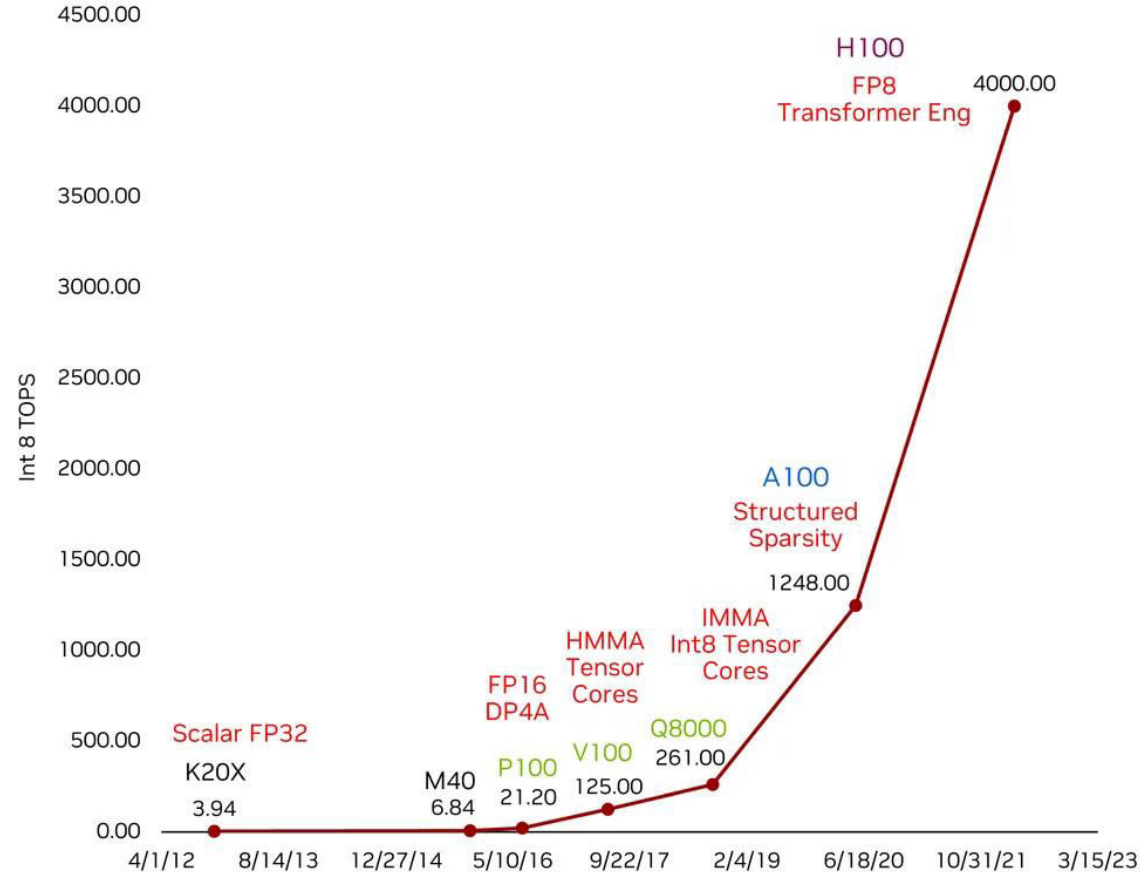  - Chiplets, 3D integration

# The Renaissance of Design

Dally HotChips 2023

Gains from

- Number Representation
  - FP32, FP16, Int8
  - (TF32, BF16)
  - ~16x

- Complex Instructions
  - DP4, HMMA, IMMA
  - ~12.5x

- Process
  - 28nm, 16nm, 7nm, 5nm
  - ~2.5x

- Sparsity
  - ~2x

- Model efficiency has also improved – overall gain > 1000x



Single-Chip Inference Performance - 1000X in 10 years

ETH zürich · ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA · 4

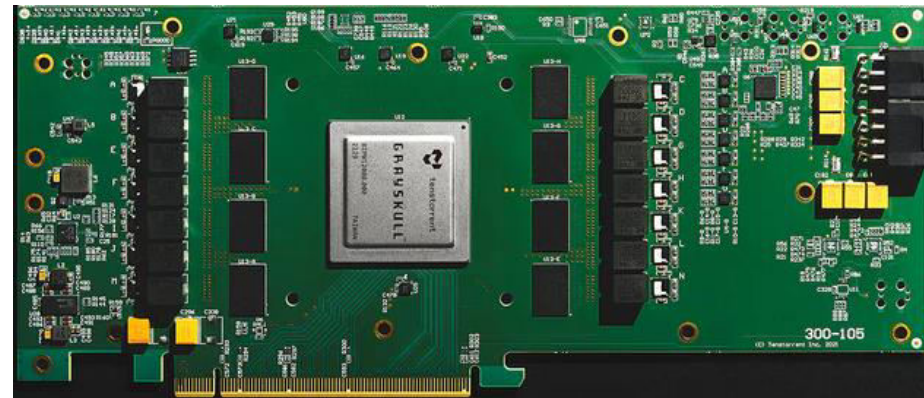# AI Innovation beyond "NVIDIA Gravity" is Challenging!

- **It's the software → flexibility, fast evolution!**

- **Need an open standard to counter a monopoly**



RISC-V: The Free and Open RISC Instruction Set Architecture
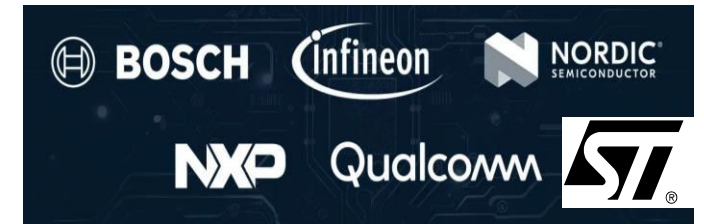
# RISC-V is Accelerating

EuroHPC 200+M€ for RV HPC (DARE FPA)
Chips (KDT) 200+M€ for RV Automotive

India Ministry for Electronics & Information Technology launched Digital India RISC-V (DIR-V) program for commercial SHAKTI & VEGA silicon.

Industry Leaders Launch RISE to Accelerate the Development of Open Source Software for RISC-V

Six chip giants to drive **RISC-V application in automotive**, enhance industry resilience

# Heterogeneous Specialization for AI

*Brain-inspired*: Multiple areas, different structure different function!



**Higher Mental Functions**
Concentration
Planning
Judgment
Emotional expression
Creativity
Inhibition - Ability to control self

**Motor Function Area**
Eye movement and placement of eyes

**Broca's Area**
Ability to talk
Ability to write

**Motor Function Area**
Ability to move muscles

**Association Area**
Short-term memory
Emotion

**Sensory Area**
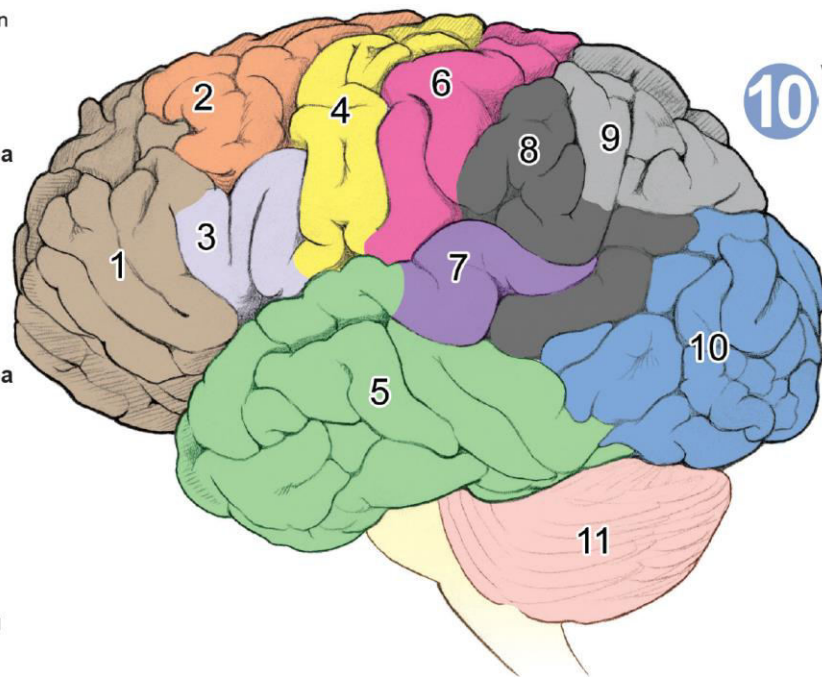Touching and feeling

**Auditory Area**
Hearing

**Wernicke's Area**
Written and spoken language understanding

**Somatosensory Association Area**
Understanding of weight, texture, temperature, etc. for recognizing and comprehending an object

**Visual Areas**
Sight
Ability to recognize pictures
Awareness of size and shape

**FUNCTIONAL AREAS OF THE CEREBELLUM**

**Motor Functions**
Coordination of movement
Balance
Posture

**Multi-sensory input**
**Frame-based**
**Event based**

➡ **Perception**
**Fusion**
**Reasoning**

ETH zürich    ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA

# Focus on Processing Element: Specialize

**RISC-V** **Instruction set: open and extensible *by construction* (great!)**

8-bit Convolution

### Vanilla

```
addi   a0,a0,1
addi   t1,t1,1
addi   t3,t3,1
addi   t4,t4,1
lbu    a7,-1(a0)
lbu    a6,-1(t4)
lbu    a5,-1(t3)
lbu    t5,-1(t1)
mul    s1,a7,a6
mul    a7,a7,a5
add    s0,s0,s1
mul    a6,a6,t5
add    t0,t0,a7
mul    a5,a5,t5
add    t2,t2,a6
add    t6,t6,a5
bne    s5,a0,1c000bc
```

**RISC-V core**

N

### Specialized for AI

```
Init NN-RF (outside of the loop)
lp.setup
pv.nnsdotup.h  s0,ax1,9
pv.nnsdotsp.b  s1, aw2, 0
pv.nnsdotsp.b  s2, aw4, 2
pv.nnsdotsp.b  s3, aw3, 4
pv.nnsdotsp.b  s4, ax1, 14
end
```

**RISC-V core**

N/4

**15x** less instructions than Vanilla!

**Specialization Cost: Power,Area: 1.5x↑ but Time 15x↓  →  E = PT 10x ↓**
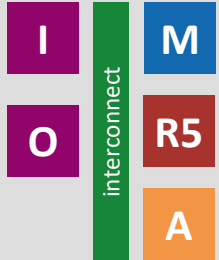
# Fully Open-Source SW Stack for AI with RISC-V!



**QuantLab**
Quantization Laboratory

**Deeploy**

**PULP-NN**
**PULP N**eural
**N**etwork backend

Specification and dataset selection

Training

Quantization/Pruning

PyTorch

ONNX

Graph optimization

Memory-aware deployment

Accelerator mapping

Tiling

Accelerator mapping

MLIR

Optimized DNN library

LLVM
COMPILER INFRASTRUCTURE

RISC-V®

# Open-Source RISC-V Hardware: PULP

## RISC-V Cores and Vector Units

| RI5CY CV32E | Zero R Ibex | Snitch | Spatz | Ariane CVA6 | ARA |
|---|---|---|---|---|---|
| RV32 | RV32 | RV32 | RVV | RV64 | RVV |

## Peripherals

| JTAG | SPI |
|---|---|
| UART | I2S |
| DMA | GPIO |

## Interconnects

| LIC | HCI |
|---|---|
| APB | FlooNoC |
| AXI4 | |

## Platforms



**Single core**
- PULPino, PULPissimo
- Cheshire

**Multi-core**
- OpenPULP
- ControlPULP

**Heterogeneous, Many-core**
- Hero, Carfield, Astral
- Occamy, Mempool

IOT → HPC

## Accelerators and ISA extensions

| XpulpNN, XpulpTNN | ITA (Transformers) | RBE, NEUREKA (QNNs) | FFT (DSP) | REDMULE (FP-Tensor) |
|---|---|---|---|---|

ETH zürich    ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA

# We make everything (we can) available openly

- **All our development is on GitHub using a _permissive_ license**

  - HDL source code, testbenches, software development kit, virtual platform

## https://github.com/pulp-platform

  - Allows anyone to use, change, and make products without restrictions.

# Open RISC-V PEs:  In-order → Superscalar → OoO



CVA6

CVA6S+

XuanTie
C910

(a) Area Efficiency [GOPS/mm²]

(b) Energy Efficiency [GOPS/W]

(c) Area-Energy Efficiency [GOPS/mm²/W]

# Heterogeneous, Multiscale Accelerated Computing

## Multiple Scales of acceleration

**Extensions to processor cores**
- Explore new extensions
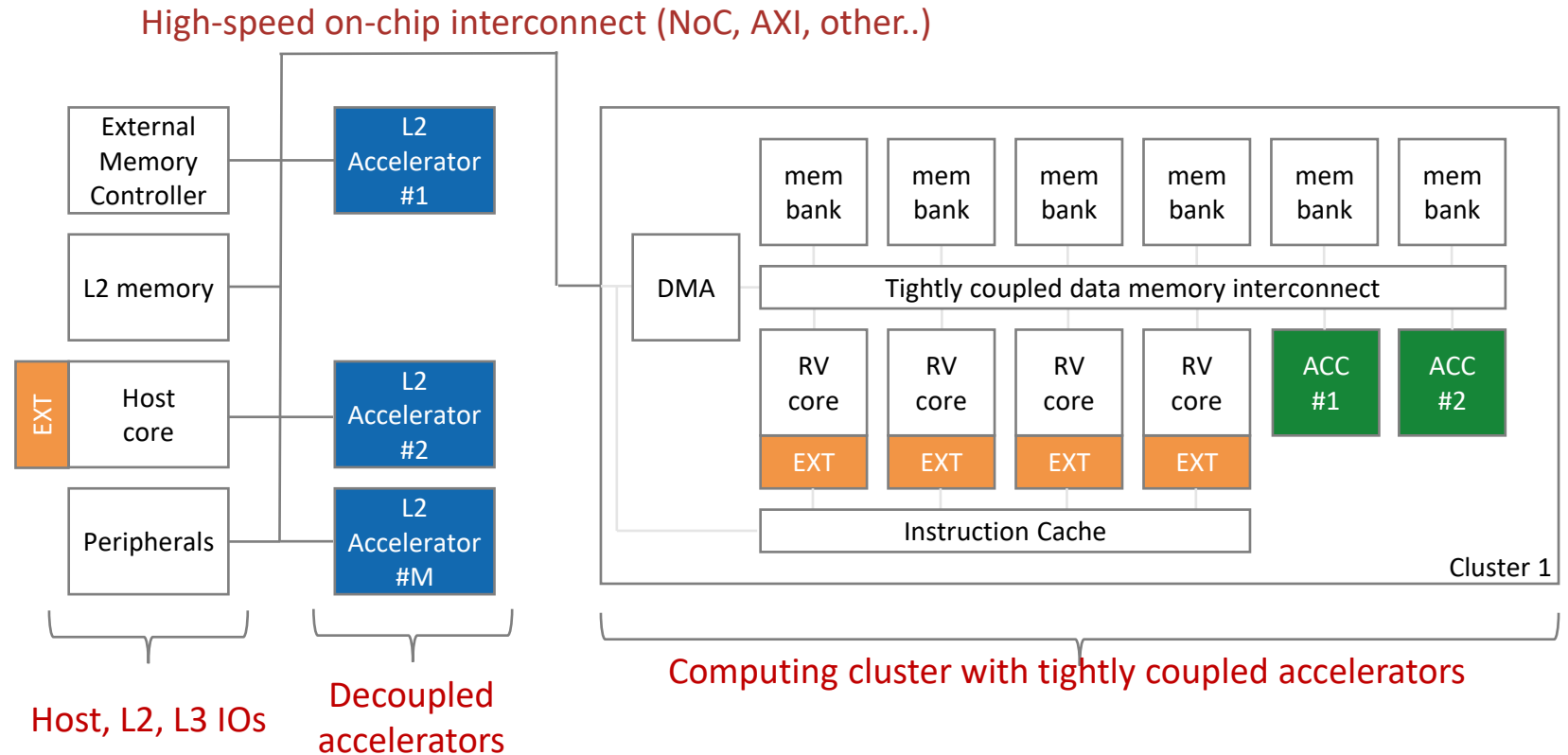- Efficient implementations

**Shared-memory Accelerators**
- Domain specific
- Local memory

**Multiple Decoupled Accelerators**
- Communication
- Synchronization

High-speed on-chip interconnect (NoC, AXI, other..)

| External Memory Controller | L2 Accelerator #1 |
| L2 memory | |
| EXT | Host core | L2 Accelerator #2 |
| Peripherals | L2 Accelerator #M |

Host, L2, L3 IOs

Decoupled accelerators

DMA

| mem bank | mem bank | mem bank | mem bank | mem bank | mem bank |

Tightly coupled data memory interconnect

| RV core | RV core | RV core | RV core | ACC #1 | ACC #2 |
| EXT | EXT | EXT | EXT | | |

Instruction Cache

Cluster 1

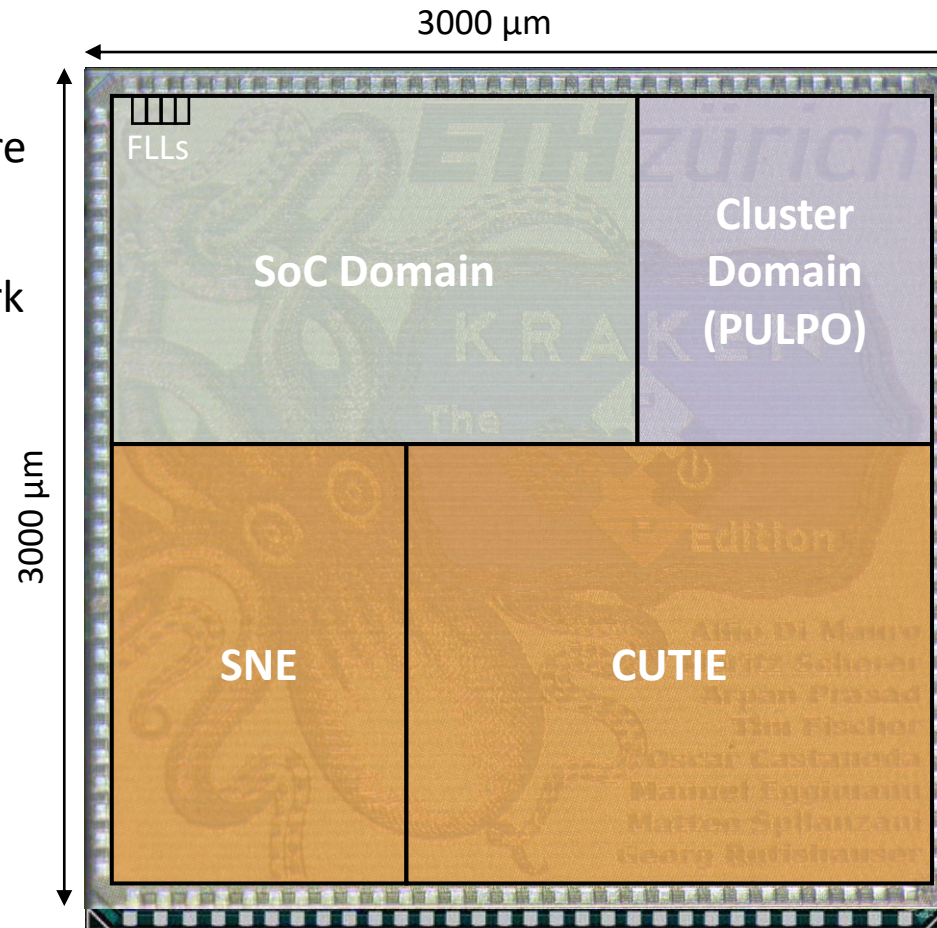Computing cluster with tightly coupled accelerators

**RISC-V is a key enabler → agility, enabling SW build-up, no vendor lock-in**

# Kraken: 22FDX SoC, Multiple Heterogeneous Accelerators

## The *Kraken*: an "Extreme Edge" Brain

- **RISC-V Cluster**

  8 Compute cores +1 DMA core

- **CUTIE**

  Dense ternary-neural-network accelerator

- **SNE**

  Energy-proportional spiking-neural-network accelerator



| Technology | 22 nm FDSOI |
|---|---|
| Chip Area | 9 mm² |
| SRAM SoC | 1 MiB |
| SRAM Cluster | 128 KiB |
| VDD range | 0.55 V - **0.8 V** |
| Cluster Freq | **~370 MHz** |
| SNE Freq | **~250 MHz** |
| CUTIE Freq | **~140 MHz** |

# Design is King

Using 22FDX tech, NT@0.6V, High utilization, minimal IO & overhead

Energy-Efficient RV Core → **10pJ (8bit)**

⬇

PE specialization 10-20x →**1pJ (8bit)**          ➡ **x10**

⬇

Configurable TE  10-20x → **100fJ (4bit)**        ➡ **x10 (x100)**

⬇

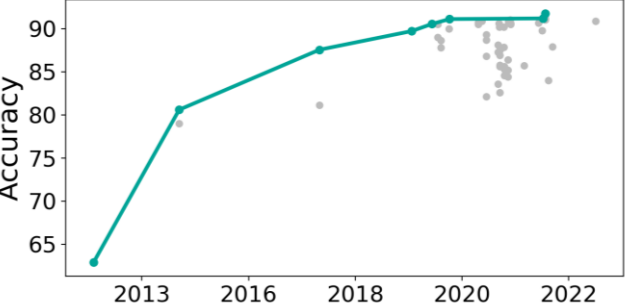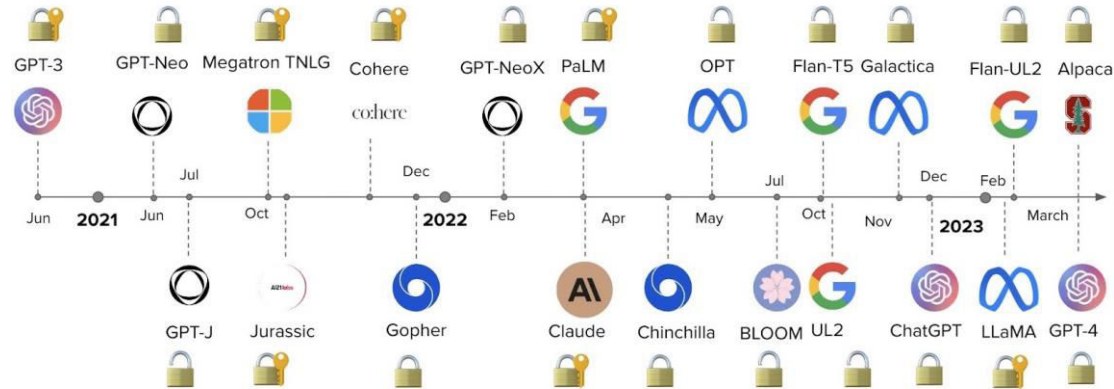Perception Accelerator 10-20x →**10fJ (events)**   ➡ **x10 (x1000)**

# Perceptive → Generative → Embodied AI



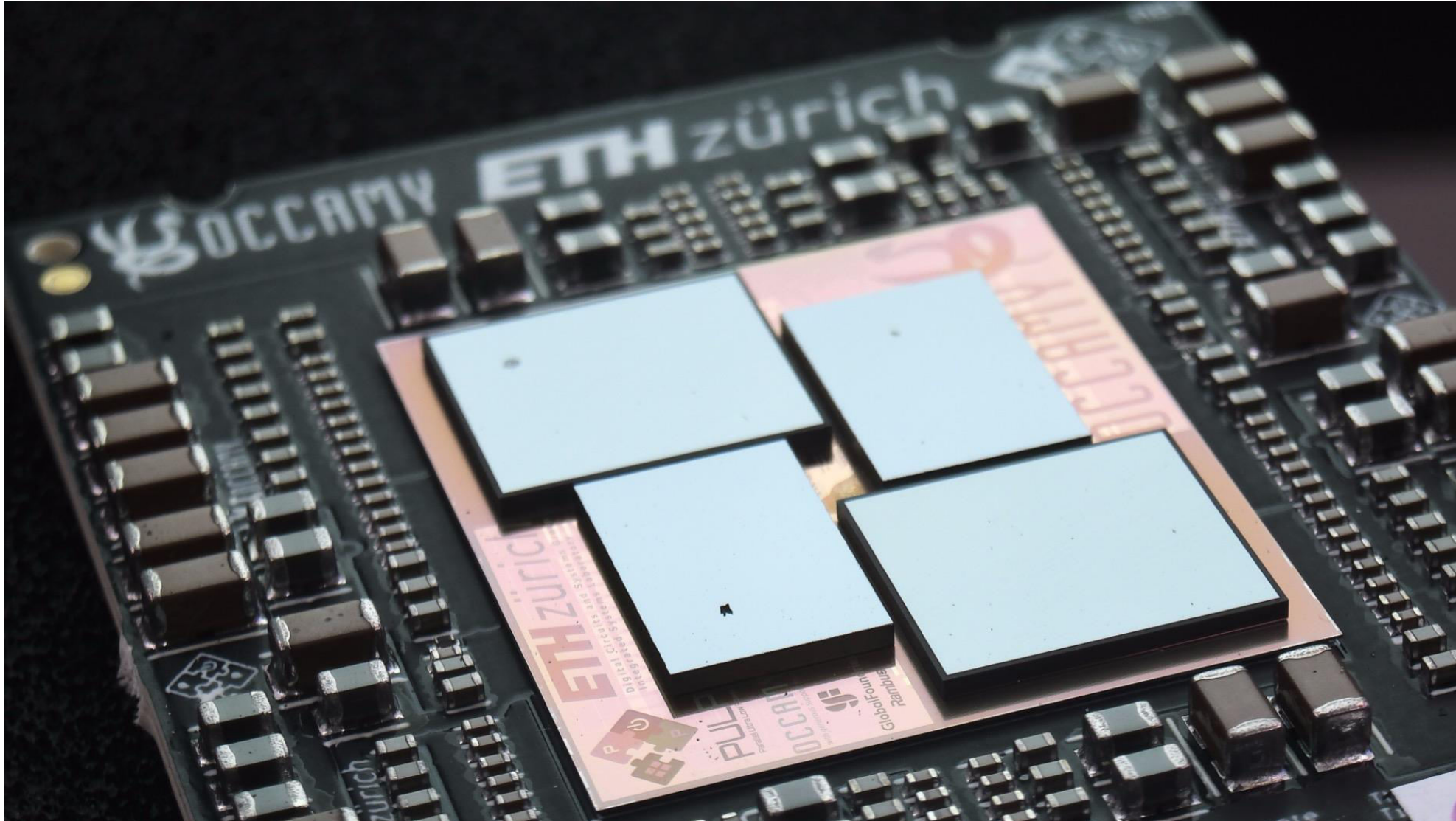Image Classification on ImageNet ReaL

**Precise**

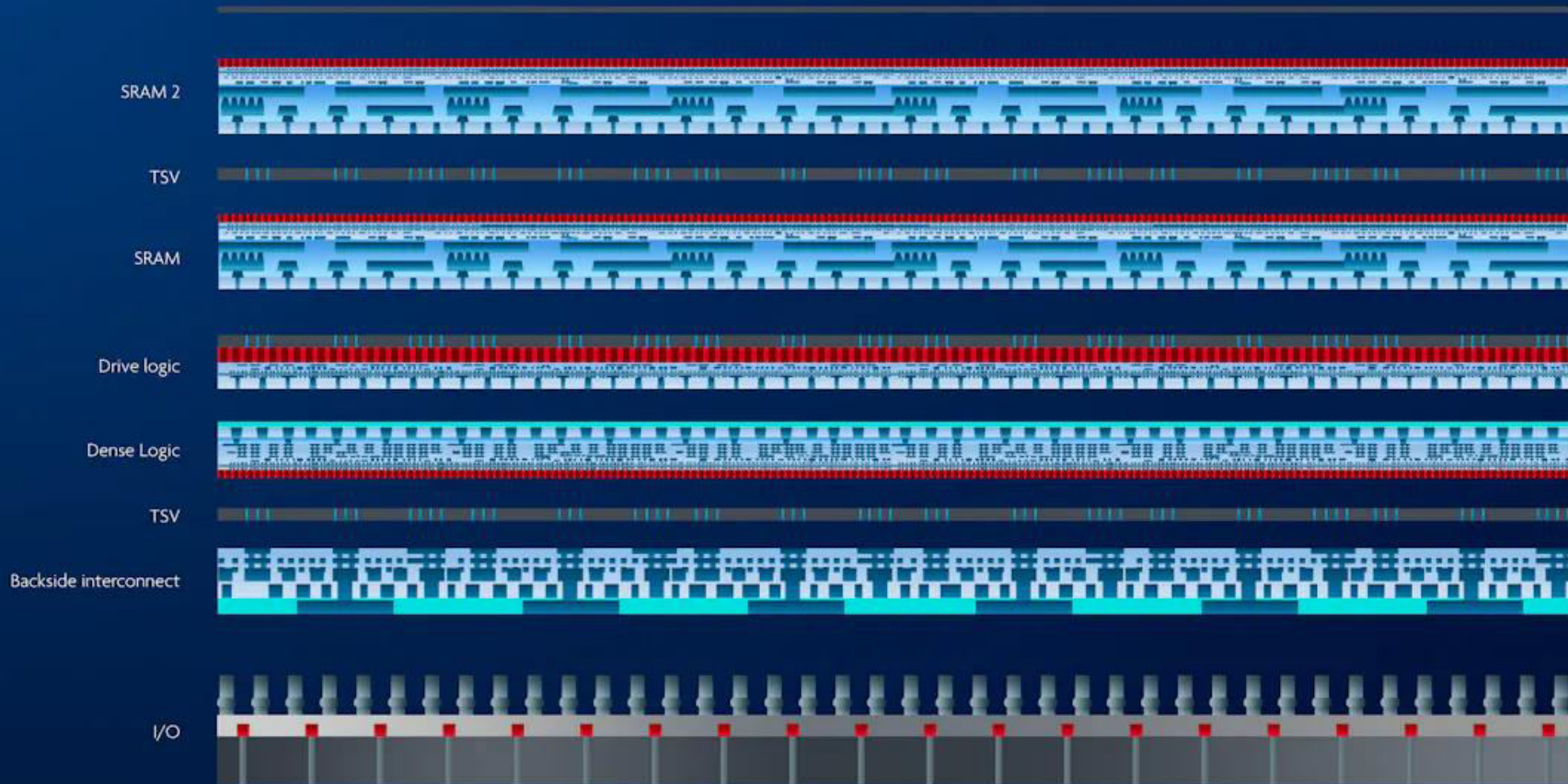**Interactive, creative**

**Efficient,
RT-safe,
secure**

# Chiplet Scaling: RISC-V gen.AI Platform for Agents

# What's next?

# What's next?  CMOS 2.0!

# Thank You!

# The Race is On: Chips Acts all around us



| USA | EU | China | South Korea | Japan | Taiwan |
|---|---|---|---|---|---|
| CHIPS for America Act | European Chips Act and IPCEI | Big Fund III | Microelectronics Cluster | Investment Fund Microelectr. industry | Taiwan Chips Act |
| $52 bn by 2026 | >$19 bn by 2030 | $41 bn[1] by 2028 | ~$3 bn[2] by 2032 | $6.8 bn by 2026 | Tax Incentives until 2029 |

Source: Strategy & analysis (LinkedIn)

**Computing for AI is the key driver!**

ETH zürich    ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA