# The **next** circuits for a better life
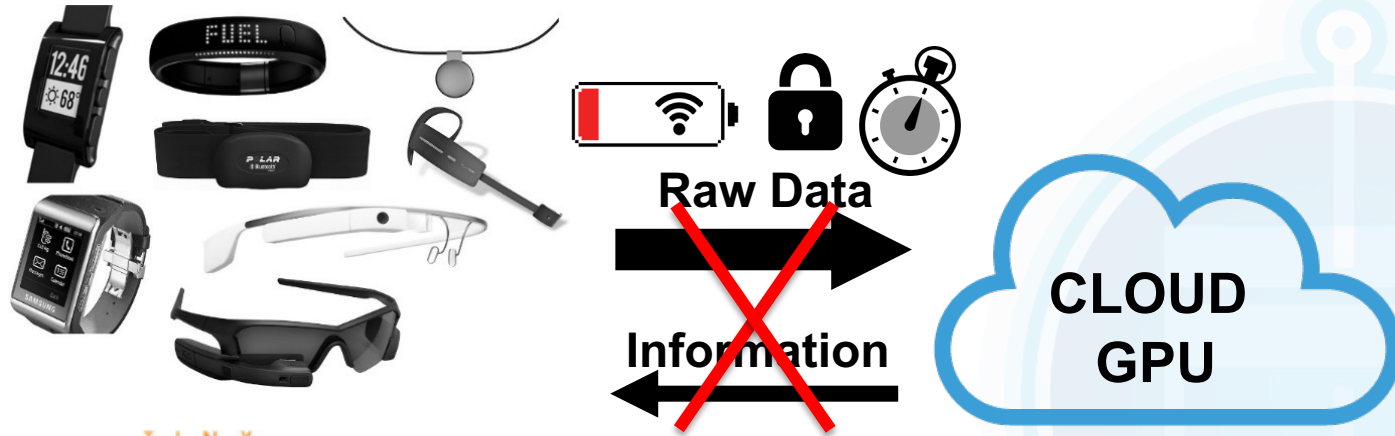
# Should tinyML Processors be Multi-core?

Marian Verhelst (marian.verhelst@kuleuven.be)

KU LEUVEN

imec

micas

# Making extreme edge (ExE) devices smart…

ExE systems = wearables, implantables, smart speaker, drones, cars, …



Raw Data

CLOUD GPU

Information

www.tinyml.org

Embedded machine learning
at the extreme edge

KU LEUVEN

micas

# Deep neural networks are everywhere in our edge devices… Are they?

Only simple tasks

KWS in phone
speech processing in cloud

Processing limited by
affordable cooling (10Watt)

Limited processing or
bulky battery

L. Lane Radius: 4.74km
R. Lane Radius: 1.09km
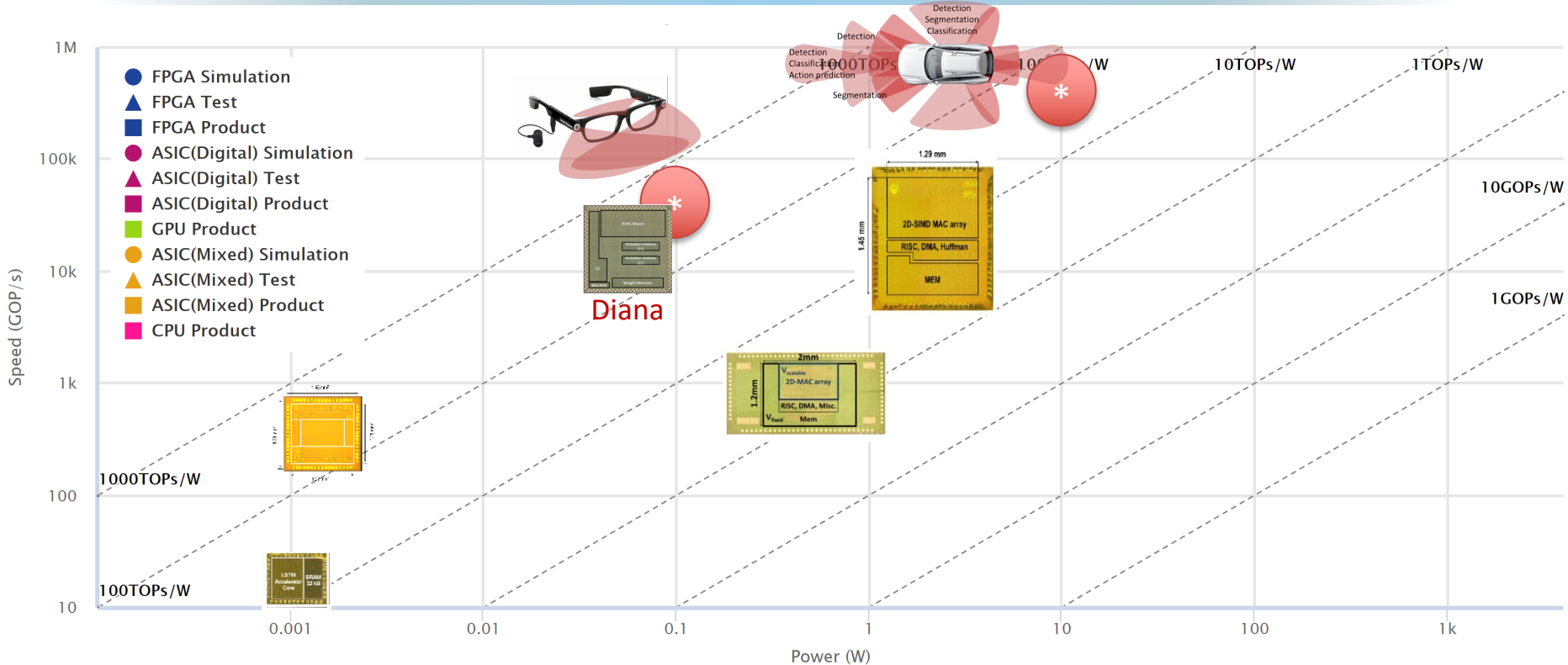C. Position: −0.40m
Close Vehicles: 2

KU LEUVEN

micas

# Deep neural networks are everywhere in our edge devices… Are they?

6 full HD cameras @30fps
10Watt, ResNet-50/frame (under est.!)

➔ 1TOPs/frame, **300 TOPs**
➔ 30TOPs/Watt

Stereo HD + eye tracking camera @30fps
100mWatt, ResNet-50/frame (under est.!)

➔ 400GOPs/frame, 30 TOPs
➔ **300TOPs/Watt**





6

KU LEUVEN

micas

# Neural network processors: state-of-the-art

# Overview

- Peak vs workload performance
  - Dependency on precision
  - Dependency on dataflows
- Motivation for heterogeneous systems
- The Diana system
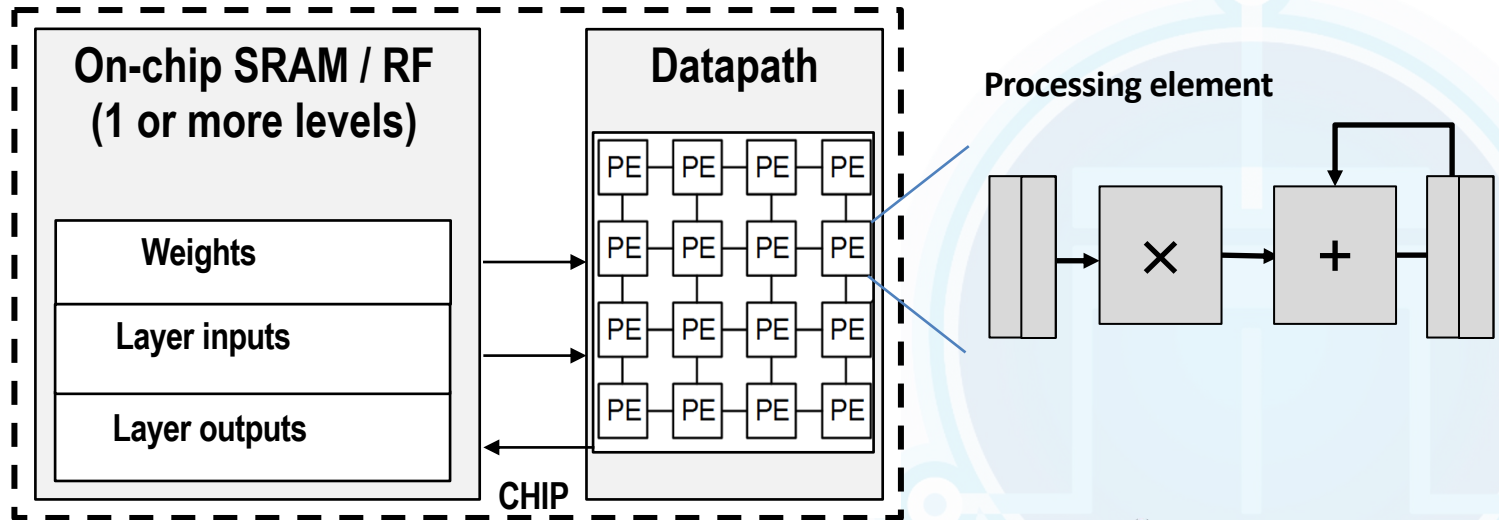- Heterogeneous scheduling with ZigZag
- The future?

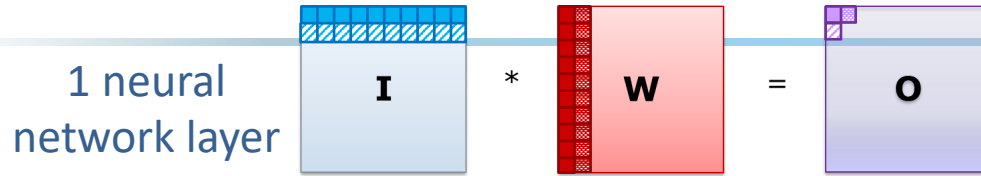# Overview

- Peak vs workload performance
  - Dependency on precision
  - Dependency on dataflows
- Motivation for heterogeneous systems
- The Diana system
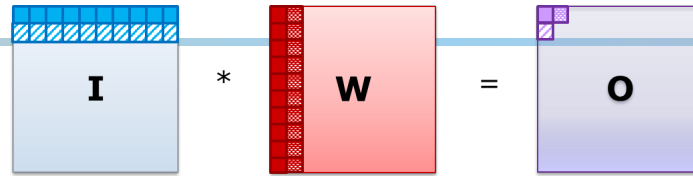- Heterogeneous scheduling with ZigZag
- The future?

KU LEUVEN

micas

# A typical Neural (co)processor unit (NPU)



1 neural network layer: $I * W = O$

On-chip SRAM / RF (1 or more levels): Weights, Layer inputs, Layer outputs

Datapath: PE array

Processing element: × +

CHIP

# A typical Neural (co)processor unit (NPU)

# "Trick" 1: Reduced precision in ML processors

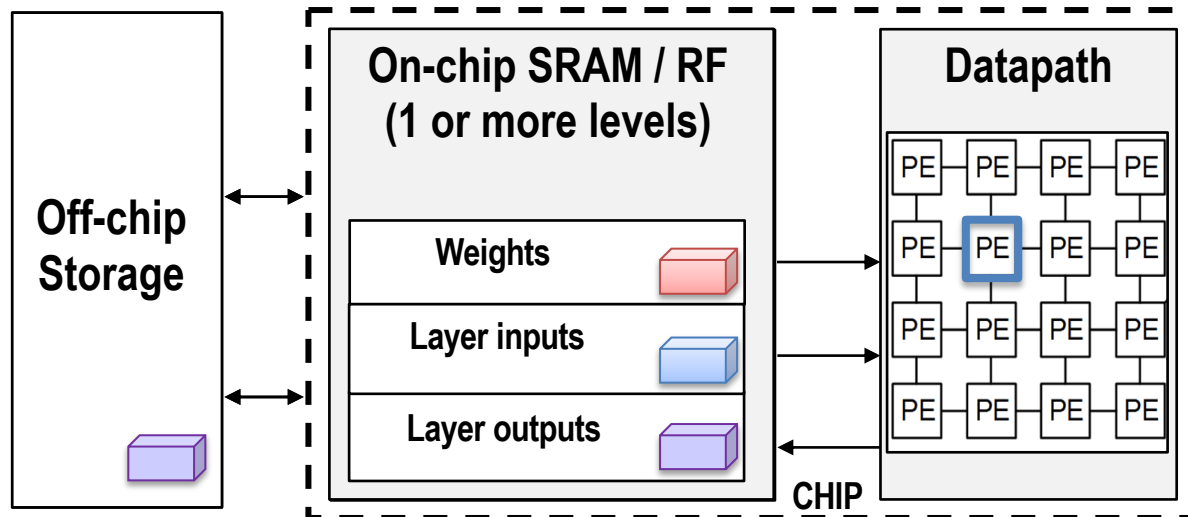I * W = O

Energy per IO transfer **/ P**

Energy per memory read/write **/ P**

Energy per MUL + ADD **/ $P^{1.5}$**

**Off-chip Storage**

**On-chip SRAM / RF (1 or more levels)**

Weights

Layer inputs

Layer outputs

**Datapath**

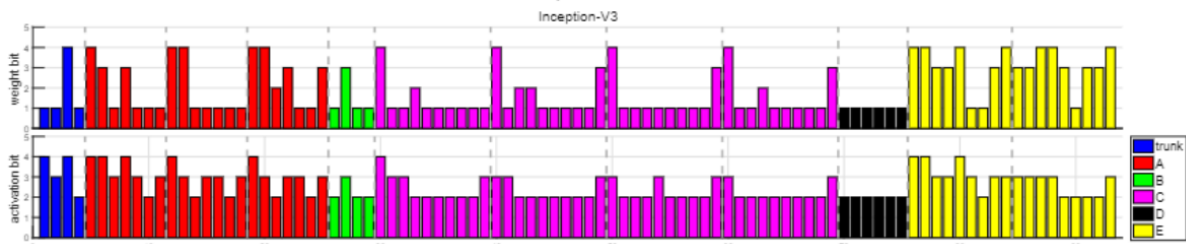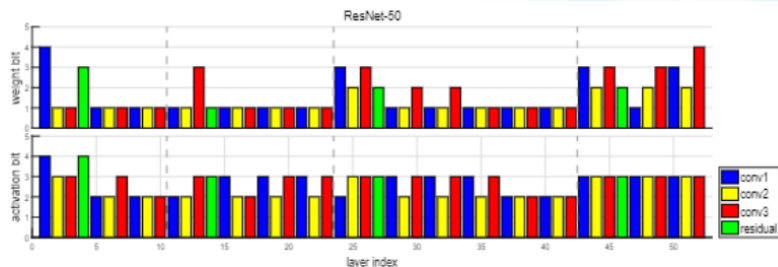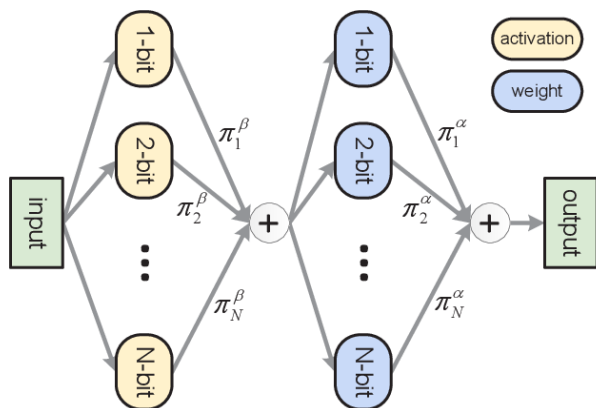| PE | PE | PE | PE |
| PE | PE | PE | PE |
| PE | PE | PE | PE |
| PE | PE | PE | PE |

CHIP

- Compute at 1-2-4 bit precision
- Exploit **precision** to reduce **memory & compute energy**
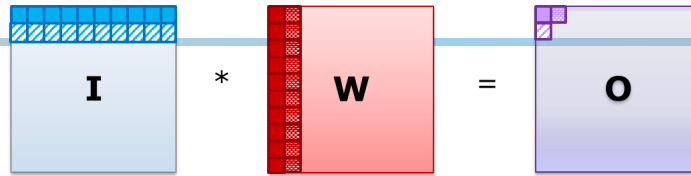
**Peak performance ~ 1/P!**

KU LEUVEN

micas

# "Trick" 1: Reduced precision in ML processors

- Active field of research in algorithmic community
- Promising results!
- Hardware support needed

### Differential architecture search



Z. Cai, N. Vasconcelos. "Rethinking Differentiable Search for Mixed-Precision Neural Networks." *CVPR,* 2020 2346-2355.

# "Trick" 2: Data reuse in ML processors



**I** * **W** = **O**

Energy per IO transfer **/N'** + Energy per memory read/write **/N** + Energy per MUL + ADD

**Off-chip Storage**

**On-chip SRAM / RF (1 or more levels)**

Weights

Layer inputs

Layer outputs

**Datapath**
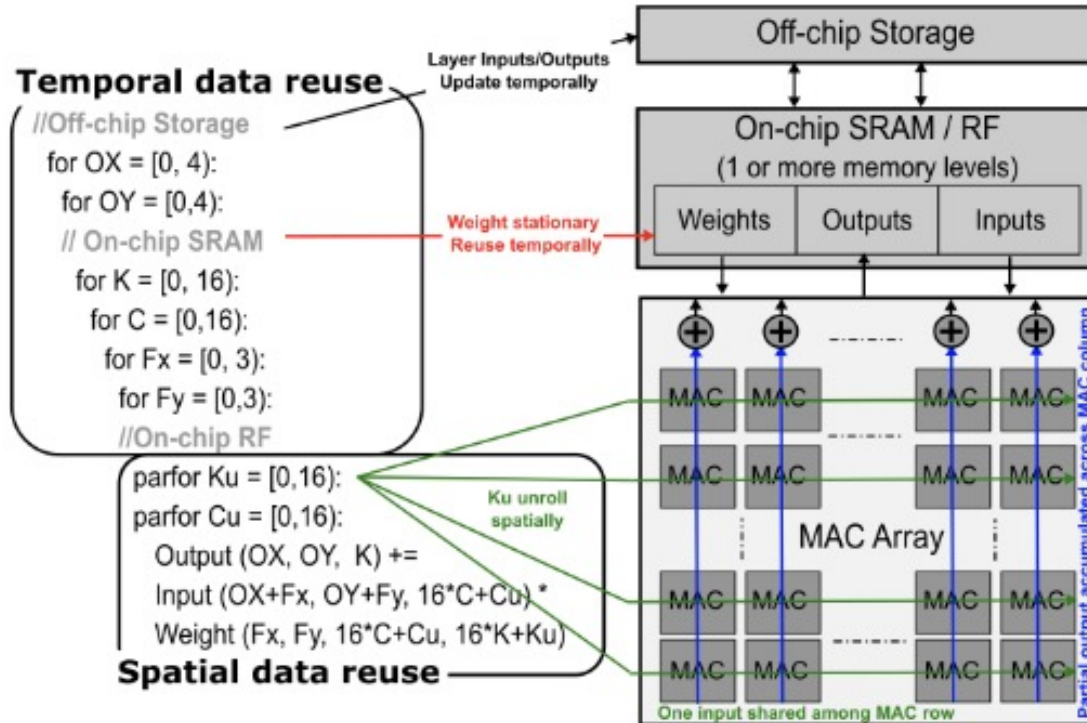
PE PE PE PE
PE PE PE PE
PE PE PE PE
PE PE PE PE

CHIP

- Remember: every W & I used multiple times, and O accumulated!
- Exploit **data reuse** to reduce **memory energy**

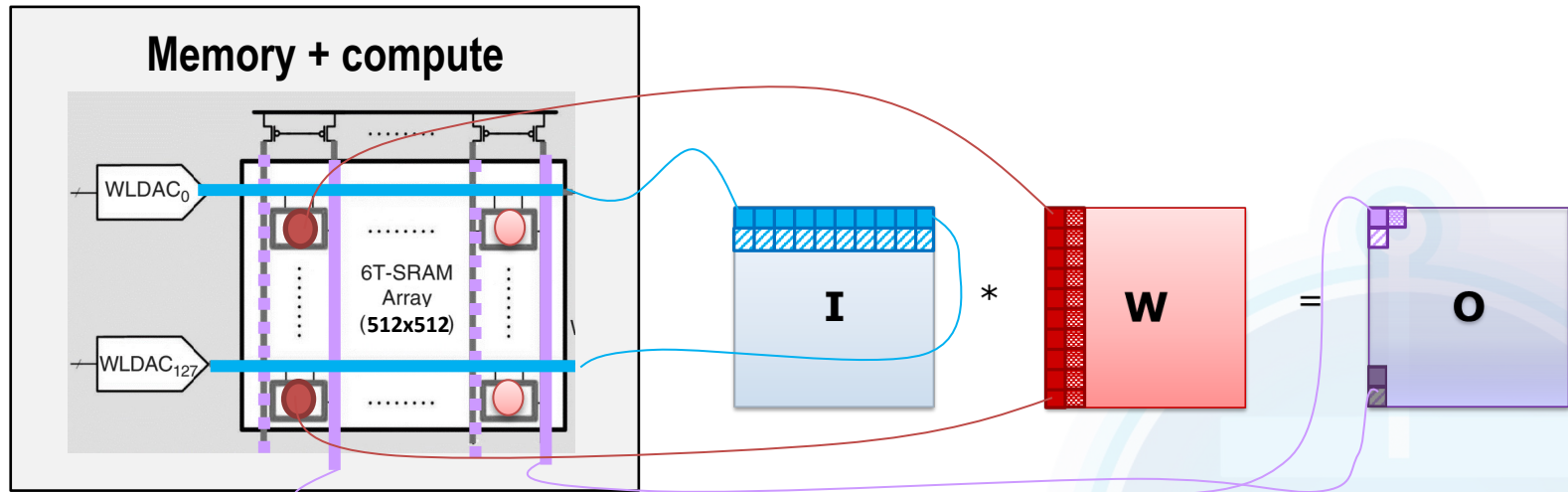**Peak performance: all data dimensions reused!**

# Data reuse in ML workloads

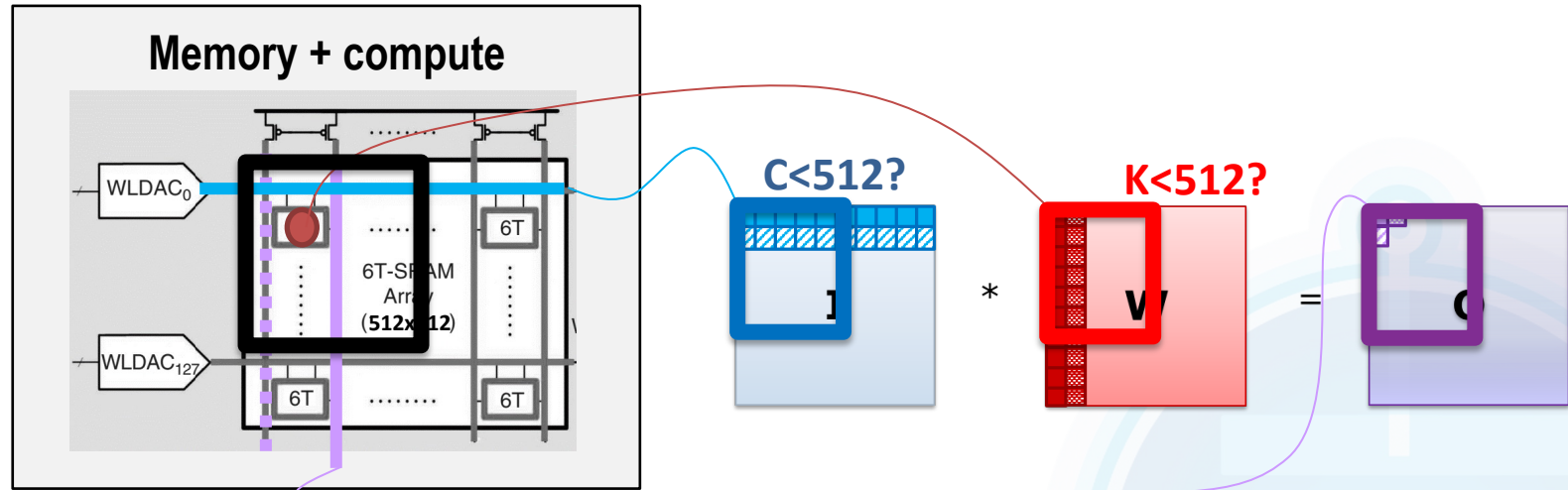Convolutions and GeMM allow significant data reuse:

# Highest peak performance? "Trick" 1 & 2!
## Analog In-memory Compute: data reuse and low precision!



- Merge the memory and compute functions
    - ➔ *bring compute to the data, instead of data to the compute!*
- Energy benefits from a.) data reuse; b.) low precision analog compute
- "Analog In-memory compute" (AiMC) ➔ **<1fJ/op**
- E.g. 512x512 size array
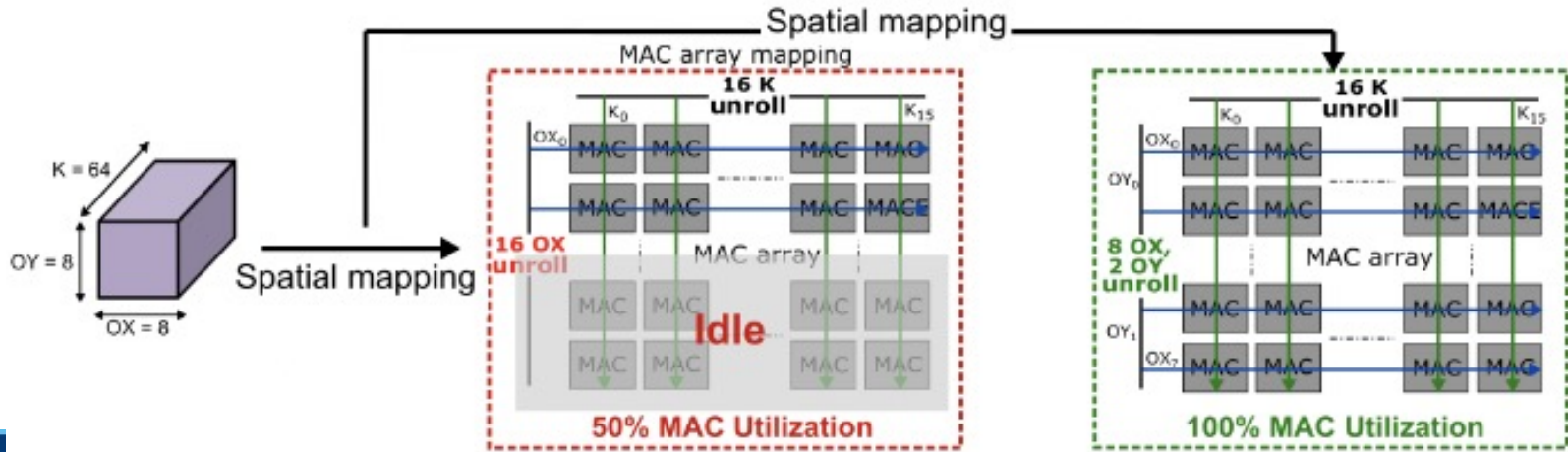
# Analog In-memory Compute: Under-utilization



- But… utilization costs!
  - Low data flow flexibility
  - Only matrices with dimensions aligned with memory array efficiently used
- E.g. 512x512 size array ➜ waste of utilization / power / …

# Efficiency for actual ML workloads

- In-memory compute enable to exploit low precision and massive parallelism
- But…:
  - But data reuse opportunities are layer dimension dependent
  - What about non GeMM layers? DW layers? FC layers? …?



Spatial Under-utilization: Workload diversity vs fixed MAC array spatial data reuse

# Neural network processors: state-of-the-art
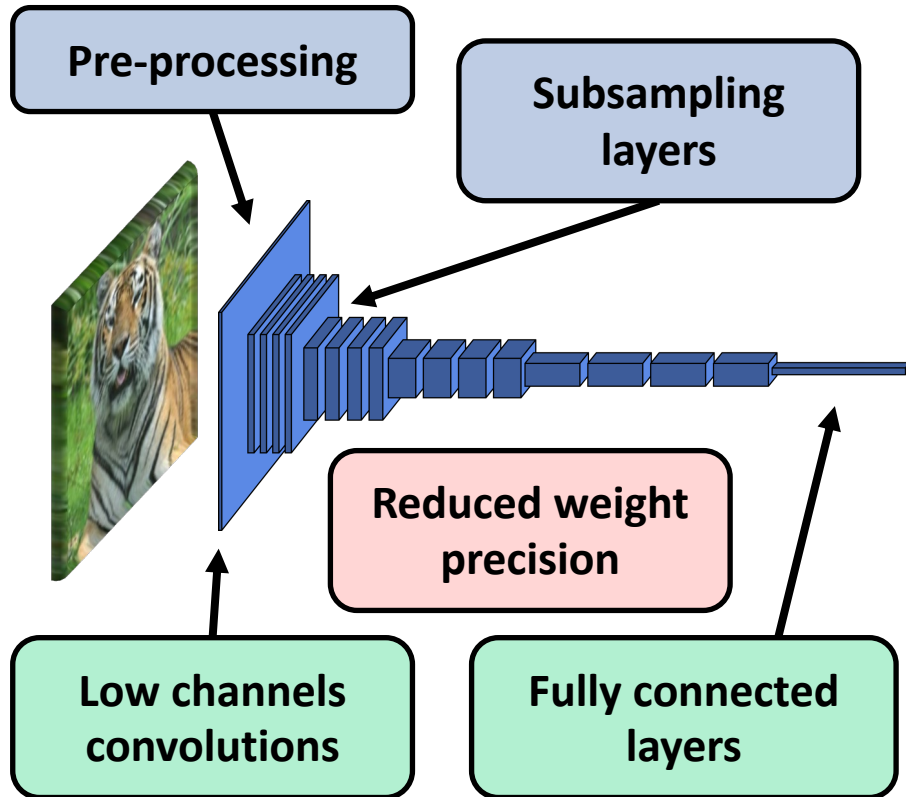
# Neural network processors: state-of-the-art



Peak performance != workload performance

# Overview

- Peak vs workload performance
  - Dependency on precision
  - Dependency on dataflows
- **Motivation for heterogeneous systems**
- The Diana system
- Heterogeneous scheduling with ZigZag
- The future?

KU LEUVEN

micas

# Need for heterogeneity



**Pre-processing**

**Subsampling layers**

**Reduced weight precision**

**Low channels convolutions**

**Fully connected layers**

**Batch norm., activation func., ...**

- **Exploit D/AiMC high energy efficiency**
  - For layers ok with **lower precision**
  - For layers with good parallelism for **high utilization**

- **BUT** have alternative accelerator(s) for other layers!

- → **Heterogeneous systems**

# Flexibility AND efficiency? ➡ Heterogeneous systems!

## Examples in the "edge"



Apple A16

**But not low power…**

Tesla FSD

# Low power for the extreme edge (tinyML)?



But lack of heterogeneity… (CPU heavy)

# Overview

- Peak vs workload performance
  - Dependency on precision
  - Dependency on dataflows
- Motivation for heterogeneous systems
- **The Diana system**
- Heterogeneous scheduling with ZigZag
- The future?

KU LEUVEN
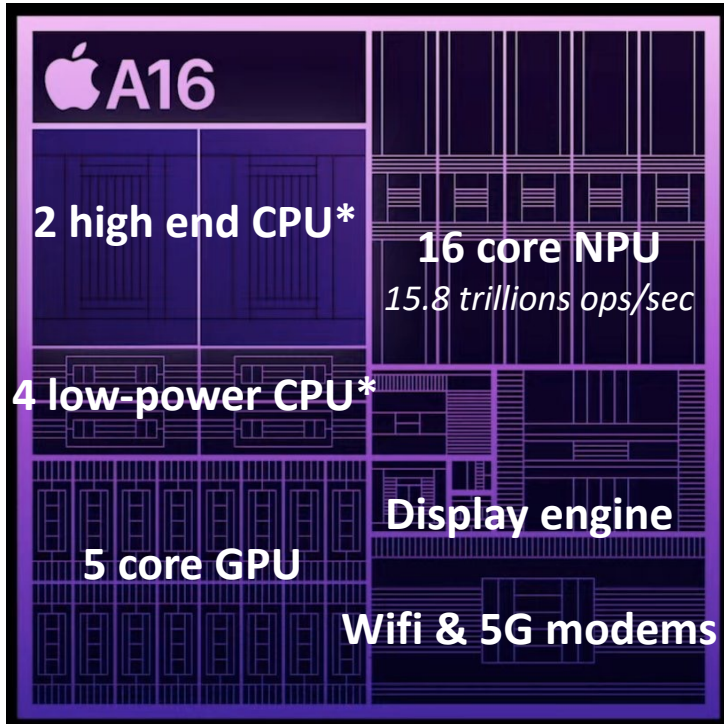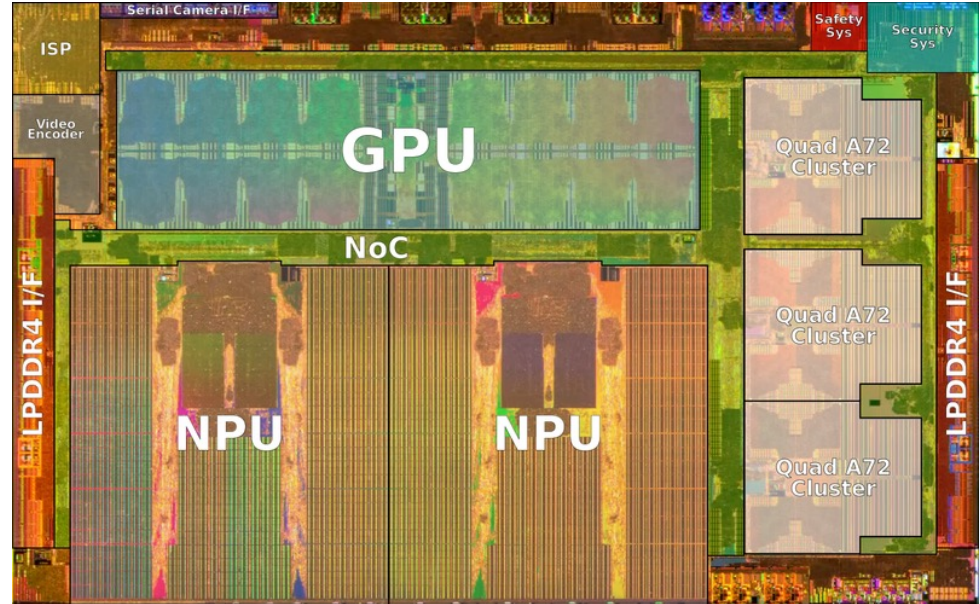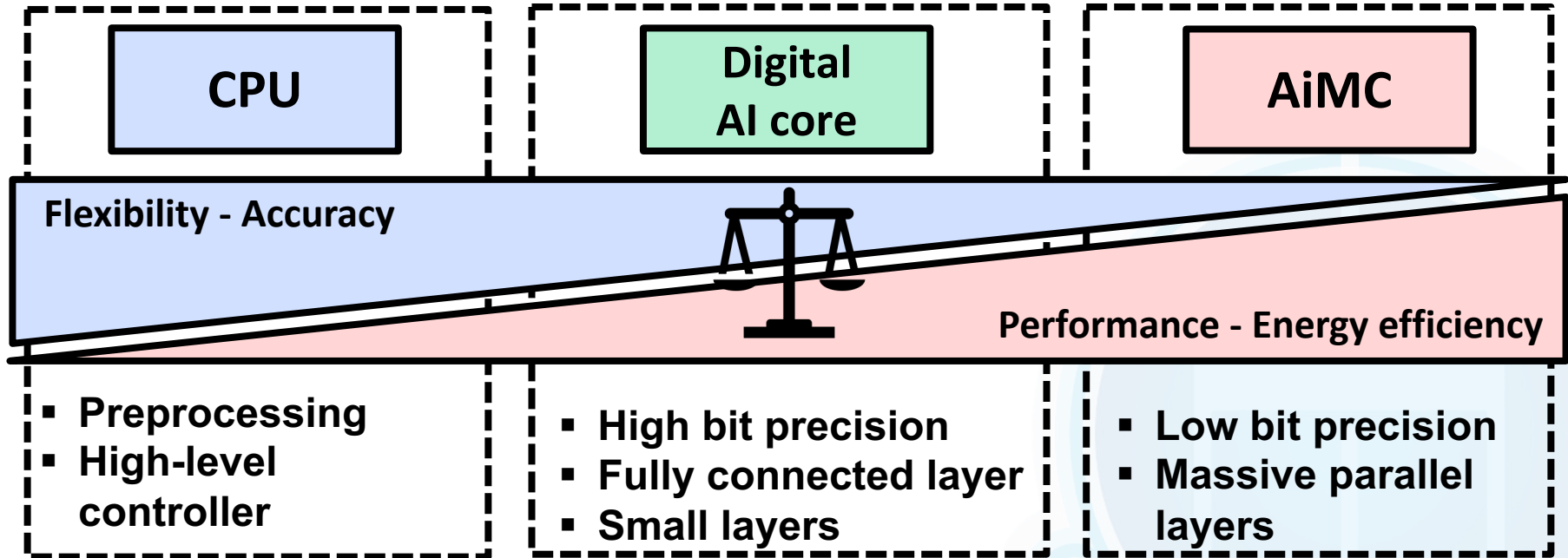
micas

# Digital-analog accelerator co-design

# DIANA SoC – High Level View

**DIANA chip:**
- RISC-V CPU*
  - High-level control
  - External I/O
- Digital AI core
  - 16x16 PEs
- Analog AI core
  - AiMC for MVMs
  - SIMD for post process
- Distributed memory hierarchy

*RISC-V CPU and periphery based on PULPissimo platform, ETH

# Digital core – Extended flexibility

**3 different levels of flexibility**

**Computation flexibility**
- 16x16 **PE array**
  - 2, 4 and 8-bit precision

**Operation flexibility**
- Convolutional layers
- Fully connected layers
- Element-wise operations
- Max pooling

**Dataflow flexibility**

# Analog core – Computing Units

**Computation units**:

- **AiMC array** SRAM-based [6] for Matrix-Vector-Multiplications
  - 1152 7-bit input DACs
  - 512 6-bit output ADCs
  - 590k compute cells (<u>ternary</u> weights)
    **Half a million MAC/cc!**

- **SIMD** for post-processing
  - 64 parallel computing units
  - 6 stages

# Analog core – Pipeline



**Three processing stages:**
1. **MCU** → Input fetch stage
2. **AiMC** → Compute stage
3. **SIMD** → Post-processing stage

**AiMC macro always in use**



**macro pipeline**, several clock cycle

# Measured results – Peak numbers

*Analog supply @ 0.8V (nominal)



**DIANA SoC - Digital core working**
8b weight - actv., 32b accumulation
**Peak efficiency/performance**

**DIANA SoC - Analog core* working**
7b DAC, 6b ADC, {-1,0,1} weights
**Peak efficiency/performance**

KU LEUVEN

# Measured results – Peak numbers

*Analog supply @ 0.8V (nominal)

**2 orders of magnitude** difference between cores
(8b digital – 2b analog weights)

| Core @ 0.8V | Efficiency (TOP/s/W) | Performance (TOP/s) |
|---|---|---|
| Digital | 2.18 | 0.177 |
| Analog | 208 | 16.9 |

DIANA So...
8b weight
**Peak eff**

Digital supply (V)

e* working
} weights
**ormance**

Nominal

Digital supply (V)

TOP/s/W

4.5
4
3.5
3
2.5
2
1.5
1
0.5
0

0.55  0.6  0.

TOP/s

20

15

10

5

0

0.85  0.9

# Overview

- Peak vs workload performance
  - Dependency on precision
  - Dependency on dataflows
- Motivation for heterogeneous systems
- The Diana system
- **Heterogeneous scheduling with ZigZag&Stream**
- The future?

**KU LEUVEN**

micas

# Hybrid execution – Pipelining

## DIANA SoC

Shared memory (**L2**)

@0.8V supplies

**CPU**
Controller
Specific operations

**Digital core**
**2.18 TOP/s/W**
**0.180 TOP/s**
High precision
Flexible

**L1**

**Analog core**
**206 TOP/s/W**
**16.85 TOP/s**
High parallelism
Limited flex.

**Different accelerators** optimized for **different workloads** to boost **system level performance**

## End-to-end mapping

Input image

layer 0

layer 1

layer n

fc

Digital core (Dc)

Analog core + SIMD

Dc

**1.) Streaming operation**
  ➔ **Efficient data sharing**
2.) Scheduling of which layer on which core?
  ➔ ZigZag!

KU LEUVEN

micas

# Hybrid execution – Layer fusion

# Hybrid execution – Layer fusion

# Layer fusion benefit on ResNet18



**ResNet18 shallow layers**

| Layer | Tile size |
|---|---|
| Conv0 (7x7) | 26x227x3 |
| Max pool. (3x3) | 12x112x64 |
| RB conv1 | 6x56x64 |
| RB conv2 | 5x56x64 |
| RB conv3 | 4x56x64 |
| RB conv4 | 4x56x64 |

**ResNet18 latency**

Latency (ms)

6.15 → 3.65  **-40%**

w/o  w

**Activation Memory Requirements**

Memory (kB)

L2

L1

980 → 195  **-80%**

w/o

- Speeds up end-to-end execution
- Reduces memory requirements

- But… scheduling degrees of freedom rise ENORMOUSLY…

KU LEUVEN

# Hybrid execution – Pipelining

## DIANA SoC

| | | |
|---|---|---|
| | **Shared memory (L2)** | |

@0.8V supplies

| **CPU** Controller Specific operations | **Digital core** 2.18 TOP/s/W 0.180 TOP/s High precision Flexible | **Analog core** 206 TOP/s/W 16.85 TOP/s High parallelism Limited flex. |

L1

**Different accelerators** optimized for **different workloads** to boost **system level performance**

## End-to-end mapping



Input image

layer 0 · layer 1 · layer n · fc

Digital core (Dc) · Analog core + SIMD · Dc

1.) Streaming operation
➔ Efficient data sharing

2.) Scheduling of which layer tile on which core at what moment?
➔ ZigZag!

KU LEUVEN

micas

# Optimizing DNN embedded processing stack with ZigZag



**ZIGZAG**

**NN workload**

**Mapping**
(spatial & temporal unrolling)

```
I   for Input
W  for Weight
O  for Output

for b = 0 to B-1        (A/O batch size)          DRAM
  for k = 0 to K-1      (O channel/ W kernel)
    for c = 0 to C-1    (A/W channel)             SRAM
      for oy = 0 to OY-1  (O row)
        for ox = 0 to OX-1  (O column)
          for fy = 0 to Fy-1  (W kernel row)      RF
            for fx = 0 to Fx-1  (W kernel column)
              unroll ...
              unroll ...                          MAC
```
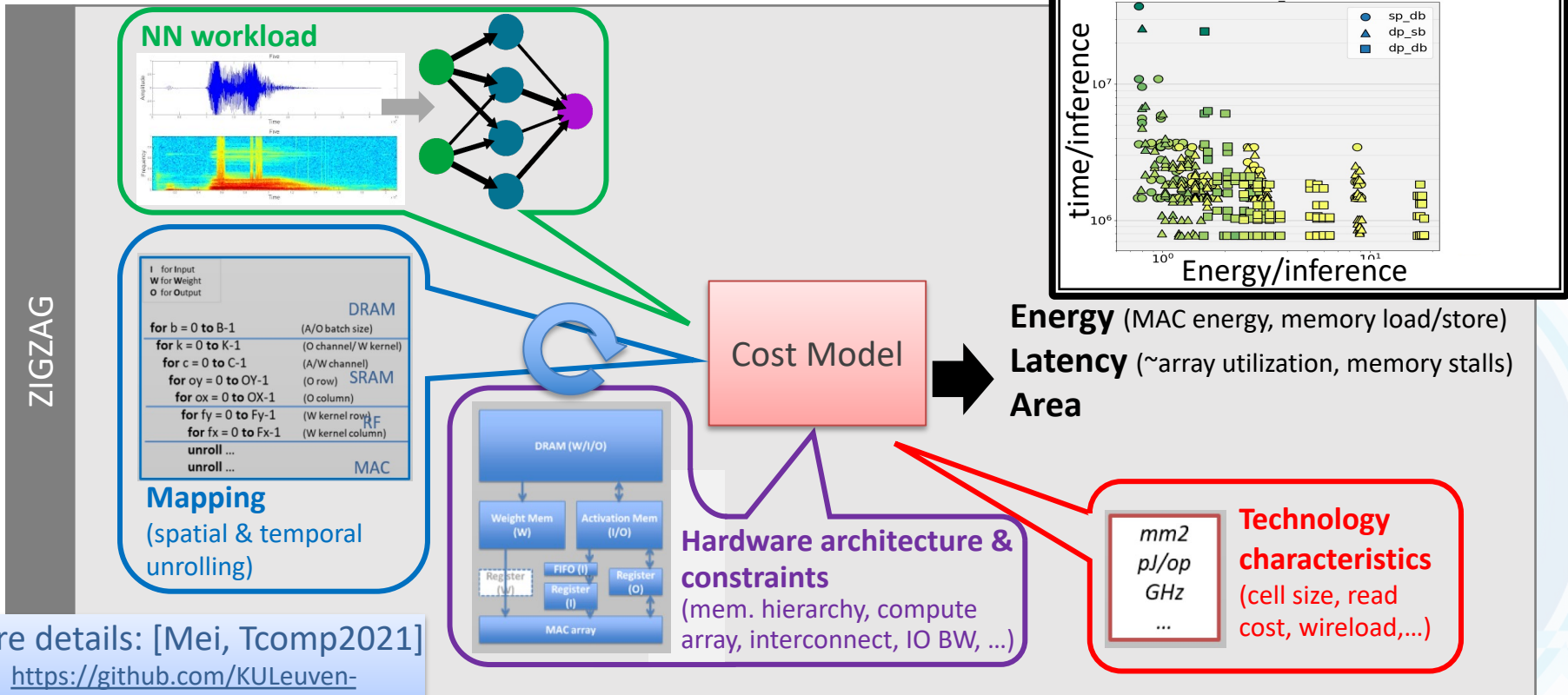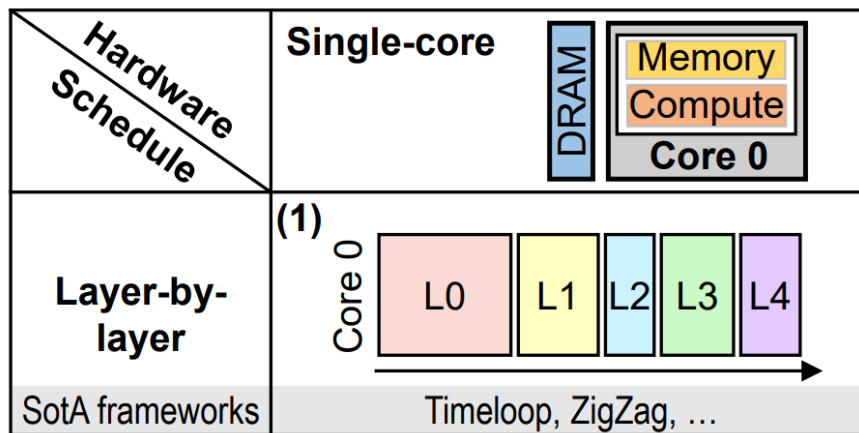
Cost Model

**Hardware architecture & constraints**
(mem. hierarchy, compute array, interconnect, IO BW, …)

DRAM (W/I/O)

Weight Mem (W)    Activation Mem (I/O)

Register (W)    FIFO (I)    Register (O)
                Register (I)

MAC array

**Energy** (MAC energy, memory load/store)
**Latency** (~array utilization, memory stalls)
**Area**

mm2
pJ/op
GHz
...

**Technology characteristics**
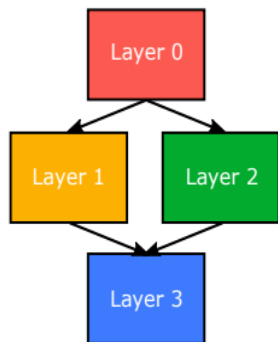(cell size, read cost, wireload,…)

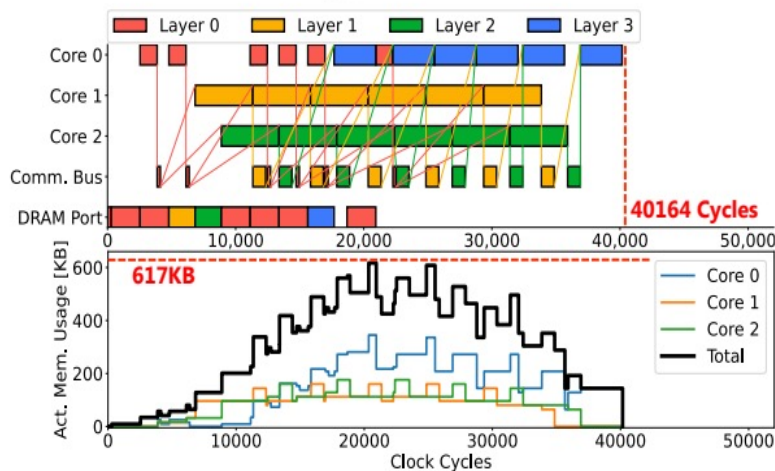# Stream: ZigZag extension to layer fusion and multi-core

# Power of ZigZag exploration:
# Scheduling optimization: latency or memory

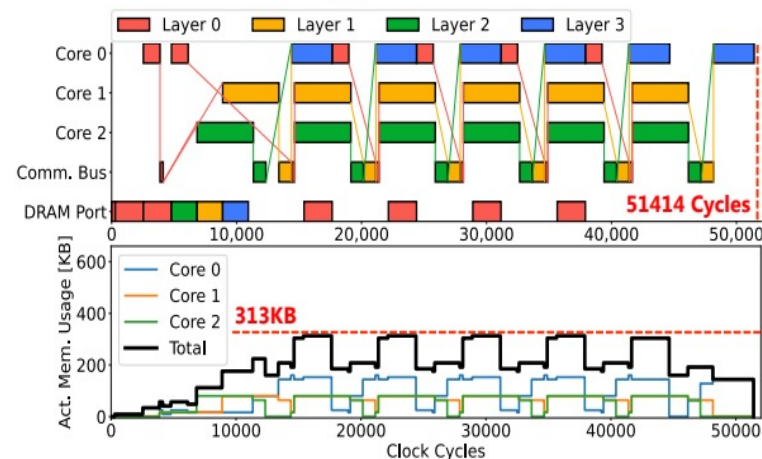- Optimize unrolling, temporal schedule, tile size, (core allocation),…



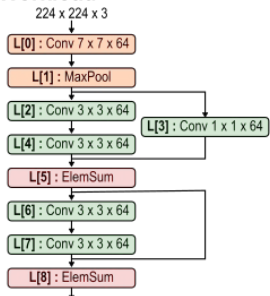(a) An example neural network

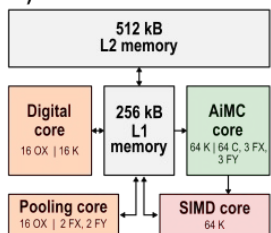(c) Latency-prioritized schedule

(d) Memory-prioritized schedule
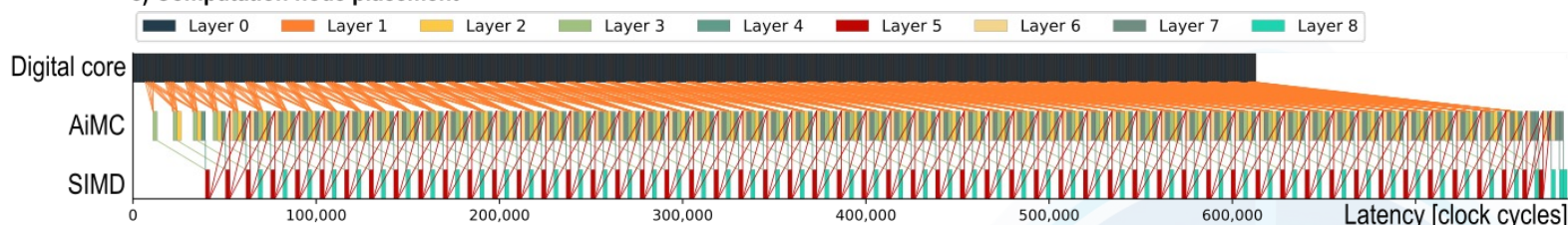
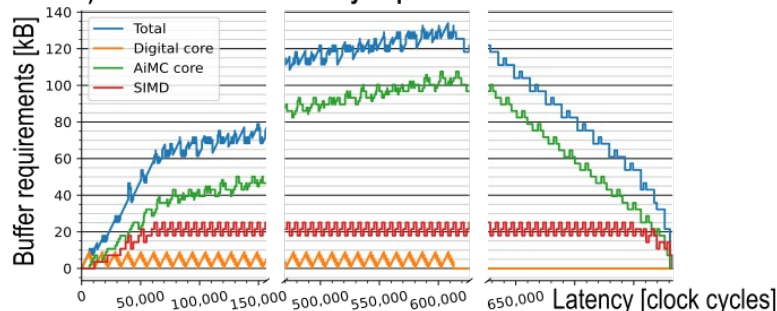# Scheduling optimization for Diana with Stream
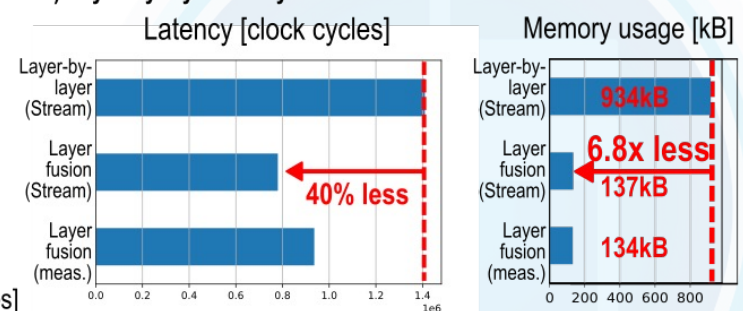


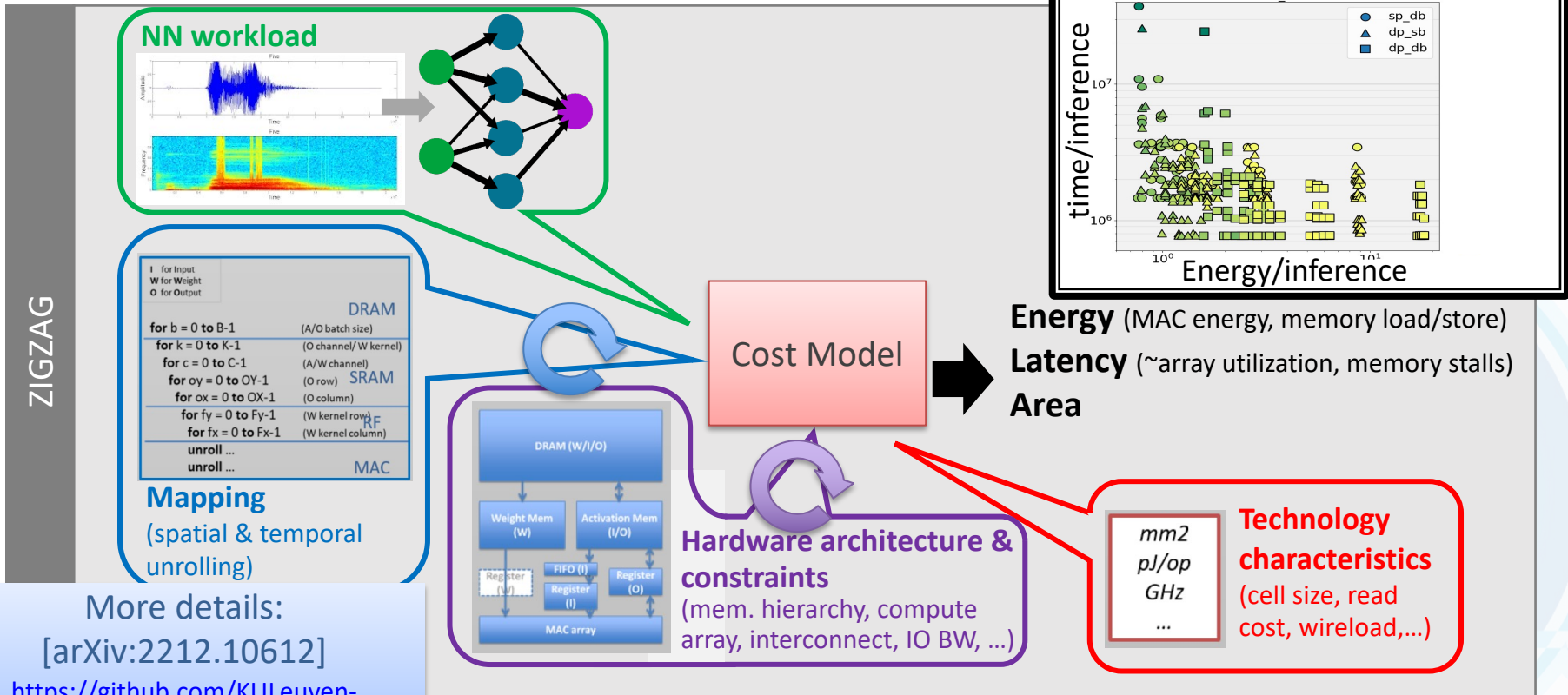a) Workload

b) Hardware architecture

c) Computation node placement

d) Activation buffer memory requirements

e) Layer-by-layer vs Layer-fusion

# Optimizing DNN embedded processing stack with Stream

**NN workload**



**Mapping**
(spatial & temporal unrolling)

```
I   for Input
W for Weight
O for Output

for b = 0 to B-1        (A/O batch size)        DRAM
  for k = 0 to K-1      (O channel/ W kernel)
    for c = 0 to C-1    (A/W channel)          SRAM
      for oy = 0 to OY-1 (O row)
        for ox = 0 to OX-1 (O column)
          for fy = 0 to Fy-1 (W kernel row)    RF
            for fx = 0 to Fx-1 (W kernel column)
              unroll ...
              unroll ...                        MAC
```

ZIGZAG

Cost Model

**Hardware architecture & constraints**
(mem. hierarchy, compute array, interconnect, IO BW, …)

DRAM (W/I/O)

Weight Mem (W)    Activation Mem (I/O)

Register (W)   FIFO (I)   Register (O)
               Register (I)

MAC array

**Technology characteristics**
(cell size, read cost, wireload,…)

mm2
pJ/op
GHz
...

**Energy** (MAC energy, memory load/store)
**Latency** (~array utilization, memory stalls)
**Area**



More details:
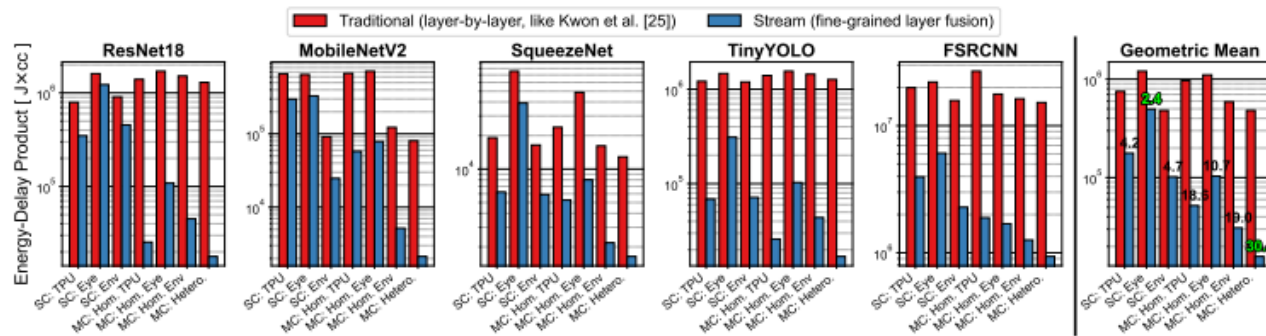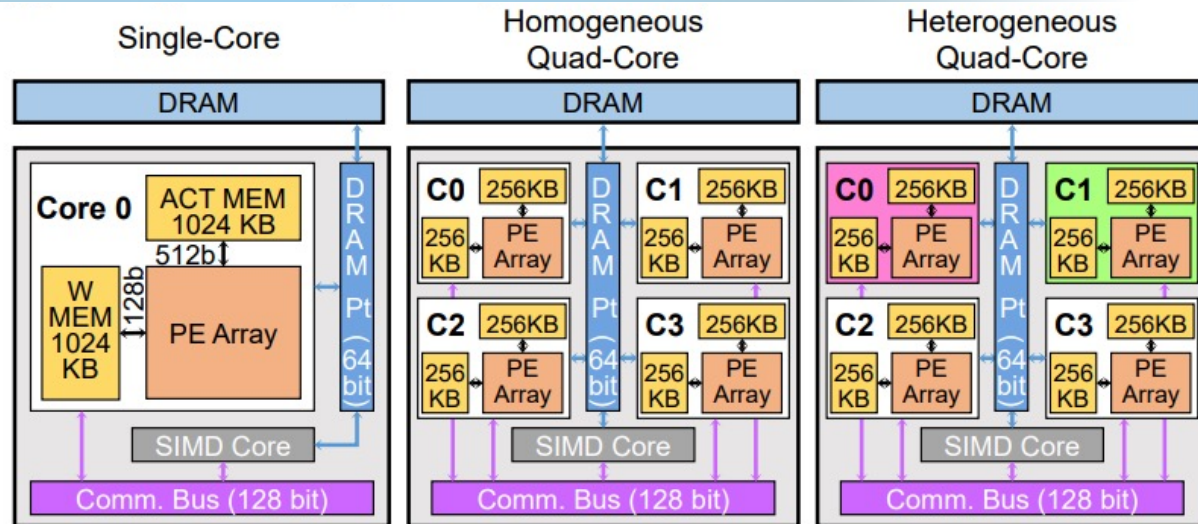[arXiv:2212.10612]
https://github.com/KULeuven-MICAS/stream

# Design space exploration

Huge design space!

- AI core sizes
- Memory volume
- Interconnect scheme
- …
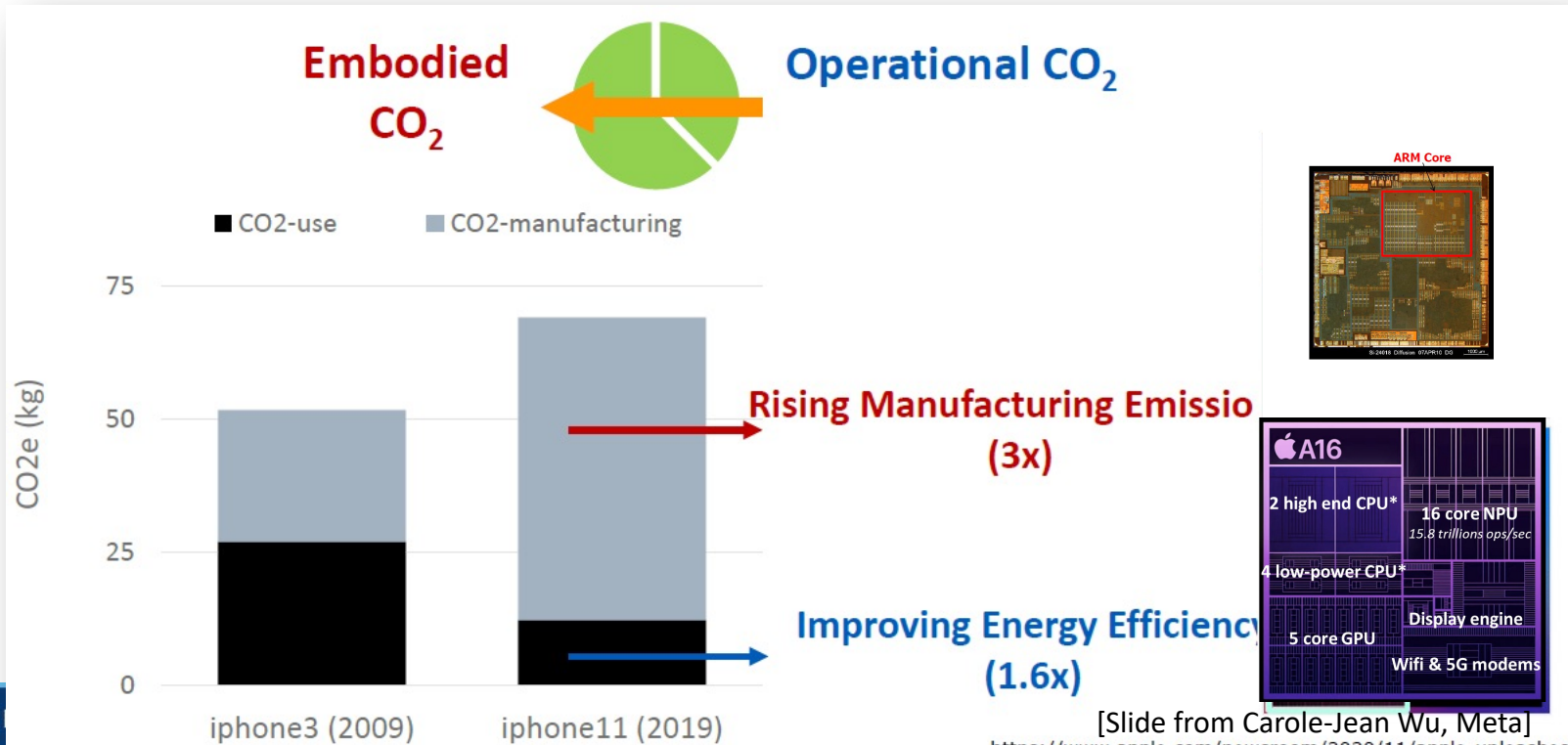
No optimal design for all networks ➡ heterogeneity helps!



KU LEUVEN

# Overview

- Peak vs workload performance
  - Dependency on precision
  - Dependency on dataflows
- Motivation for heterogeneous systems
- The Diana system
- Heterogeneous scheduling with ZigZag&Stream
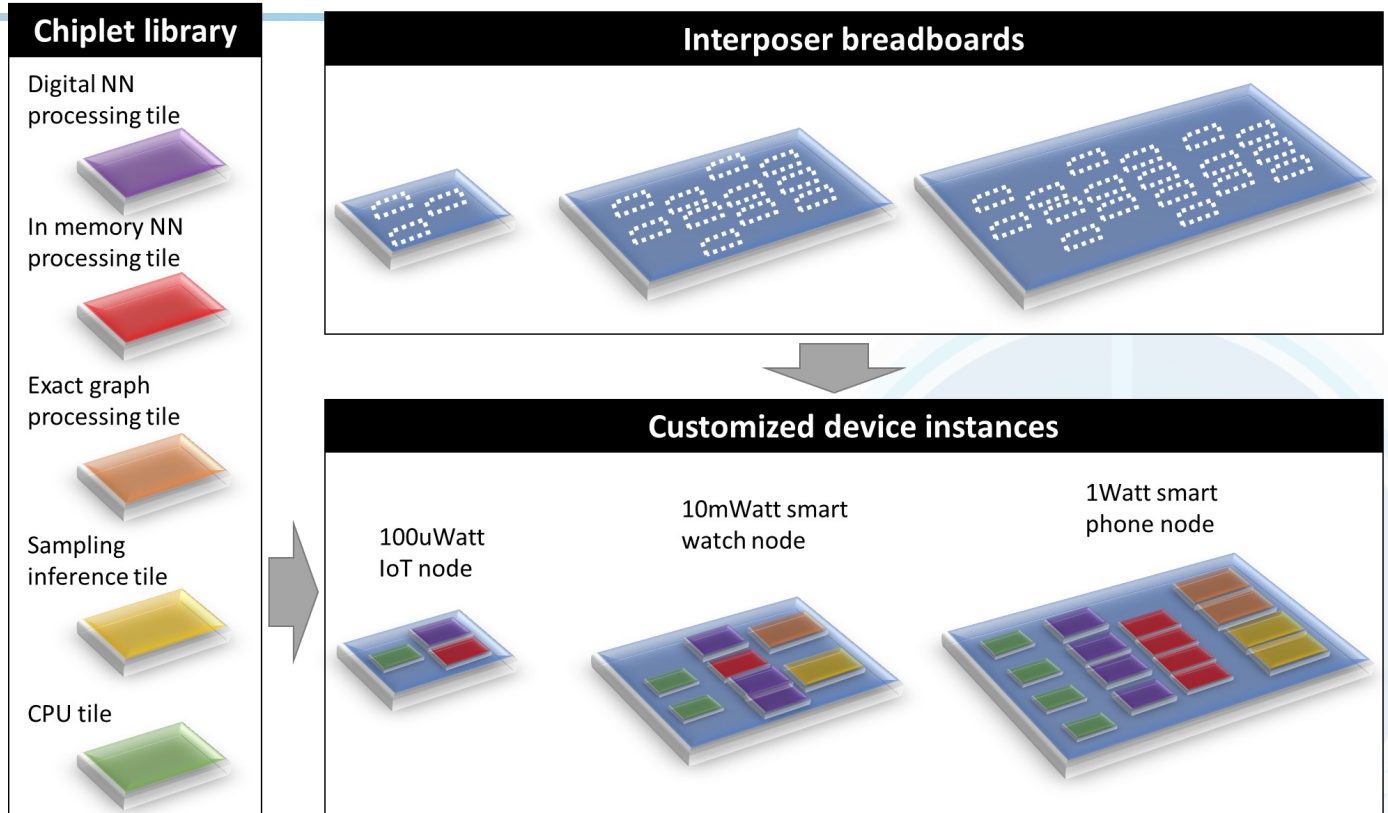- **The future?**

KU LEUVEN

micas

# Carbon footprint

- More and more accelerators on chip? No!

# Future?: More modular max-and-match with chiplets

- Supporting more workloads
- Rapid prototyping
- Customization
- Lower carbon footprint

**Chiplet library**

Digital NN processing tile

In memory NN processing tile

Exact graph processing tile

Sampling inference tile

CPU tile

**Interposer breadboards**

**Customized device instances**

100uWatt IoT node

10mWatt smart watch node

1Watt smart phone node

KU LEUVEN

micas

# Conclusion

- Specialization can bring significant efficiency gains
- **<u>Yet</u>** loss of flexibility while thriving in terms of peak performance
- ➔ Heterogeneity to have efficiency/customizability across workloads
- Towards heterogeneous, multi-core AI processing platforms
    - Rapidly customizable to algorithmic workloads
    - Supported by customizable multi-accelerator compilers
    - Large challenges ahead!

- References:
    - Ueyoshi, Kodai, et al., "DIANA: An End-to-End Energy-Efficient DIgital and ANAlog Hybrid Neural Network SoC", In 2022 IEEE International Solid-State Circuits Conference (ISSCC), IEEE, 2022.
    - Mei, Linyan, et al. "ZigZag: Enlarging joint architecture-mapping design space exploration for DNN accelerators." IEEE Transactions on Computers 70.8 (2021): 1160-1174.
    - Symons, Arneet al., "Towards Heterogeneous Multi-core Accelerators Exploiting Fine-grained Scheduling of Layer-Fused Deep Neural Networks", arXiv:2212.10612.

KU LEUVEN

micas

KU LEUVEN

imec

micas

# Copyright Notice

# www.tinyml.org