# Frontiers in Edge AI with RISC-V: Hyperdimensional Computing vs. Quantized Neural Networks
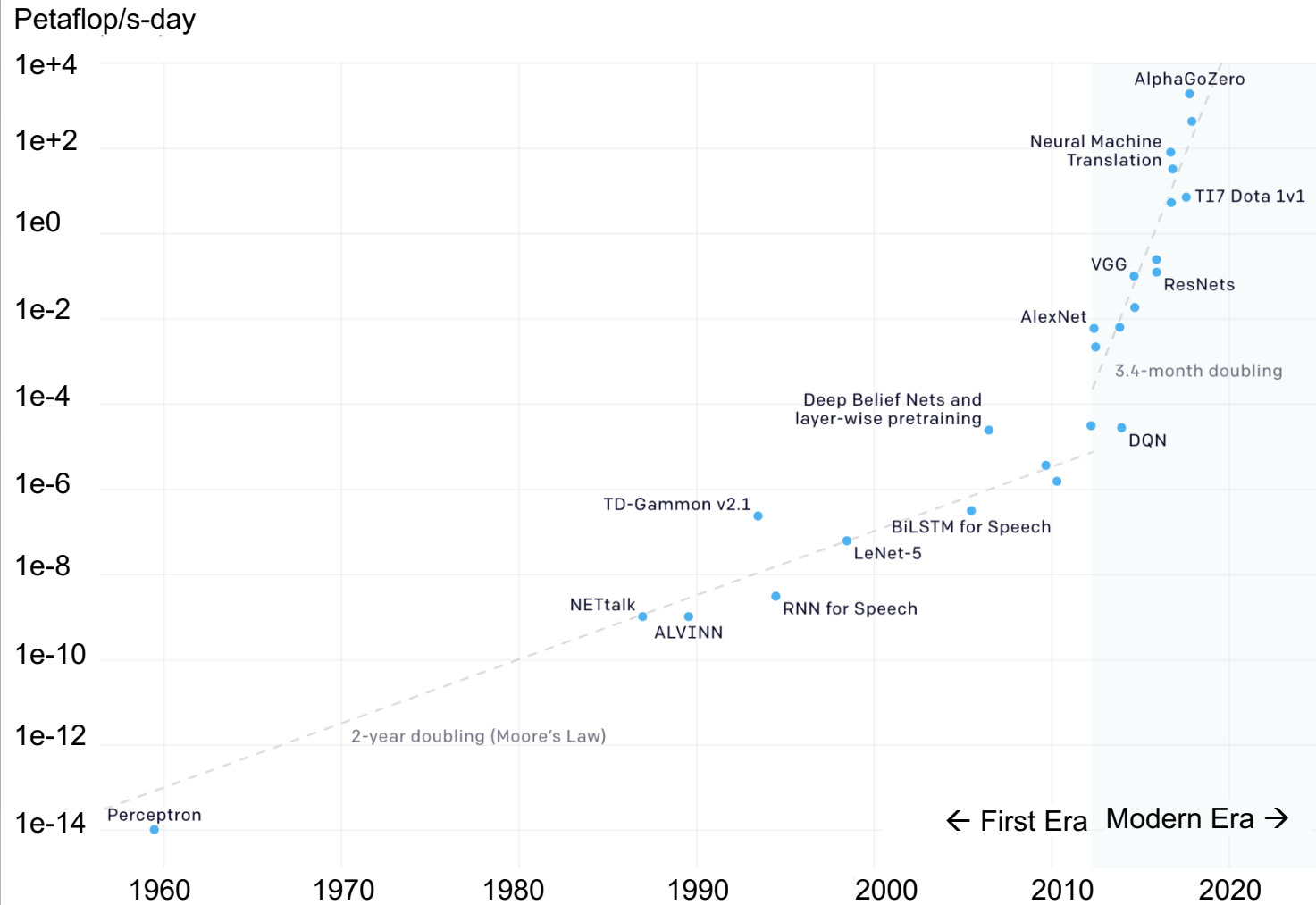
## by Hussam Amrouch
## Chair of AI Processor Design

Technical University of Munich

# The Next Revolution: AI



Source: https://openai.com

# The Next Revolution: AI

Deep Learning

by GDJ, openclipart.org

Local Unified Buffer for Activations
(96Kx256x8b = 24 MiB)
*29% of chip*

Matrix Multiply Unit
(256x256x8b=64K MAC)
*24%*

D R A M port ddr3 3%

Host Interf. 2%

Control 2%

Accumulators
(4Kx256x32b =4 MiB) 6%

Activation Pipeline 6%

PCIe Interface 3%

Misc. I/O 1%

D R A M port ddr3 3%

**AI Chip**: **Google TPUv1** [ISCA'17]

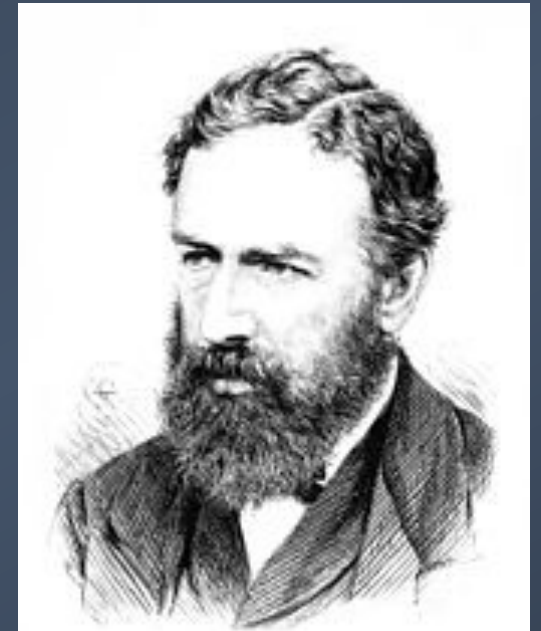**Training BERT DNN Google TPUv3: 1.8min ≈ 2048 GPUs + 512 CPUs**

# Could Efficiency be Dangerous?
# Let's go back to 1865…

## Jevons Paradox

**When technology increases the efficiency, the consumption rises.**

→ *Gain from efficiency will backfire*!

William Jevons                    src: Wikipedia

# The Upcoming Jevons Paradox

**Increase in AI Hardware Efficiency** → **Cost of DNN Training drops** → **More Companies are adopting AI**

**2030: 13% of Total $CO_2$**


src: www.pickaweb.co.uk

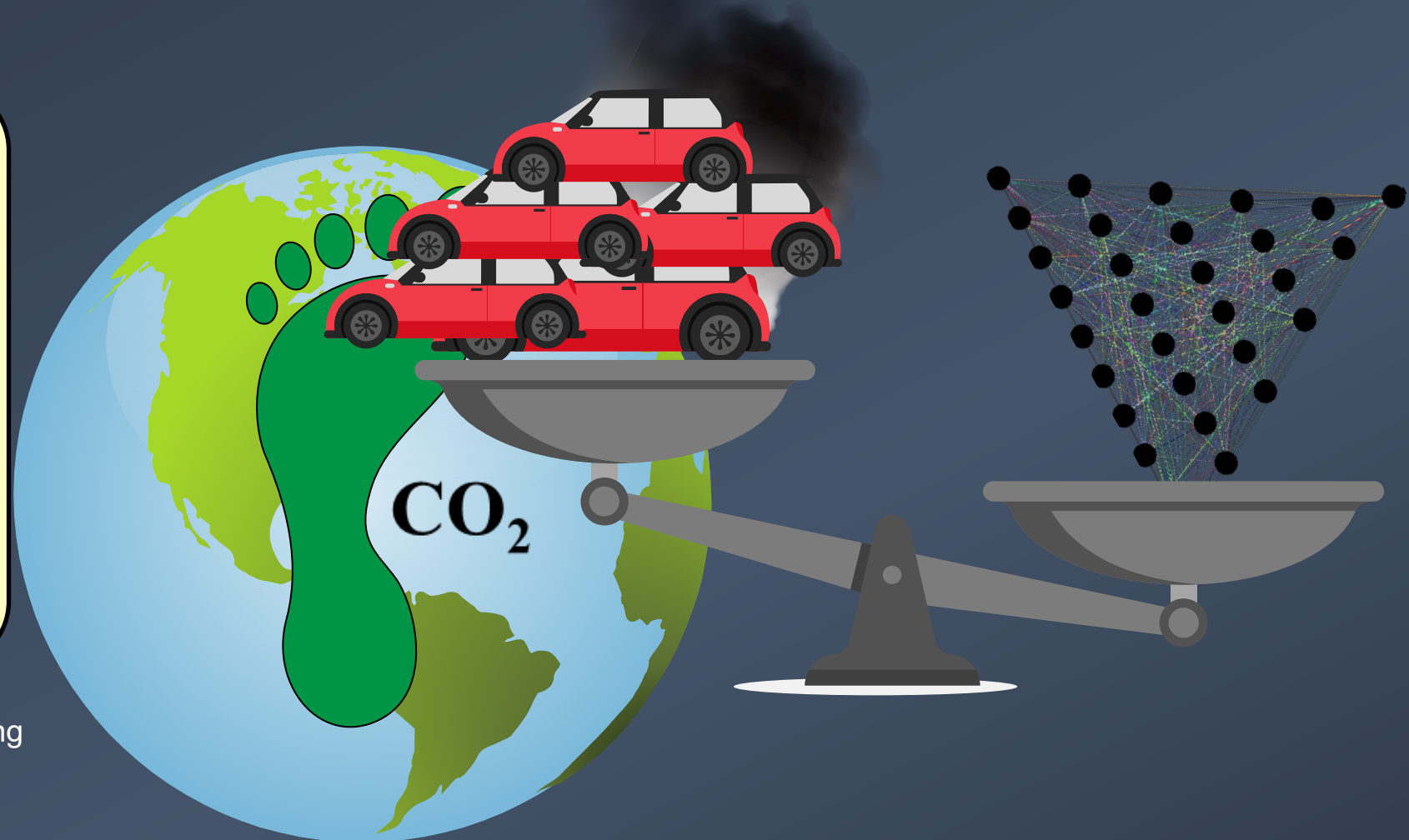**More and more data centers**

**sources: IEEE Spectrum (2019), Nature (2020)**

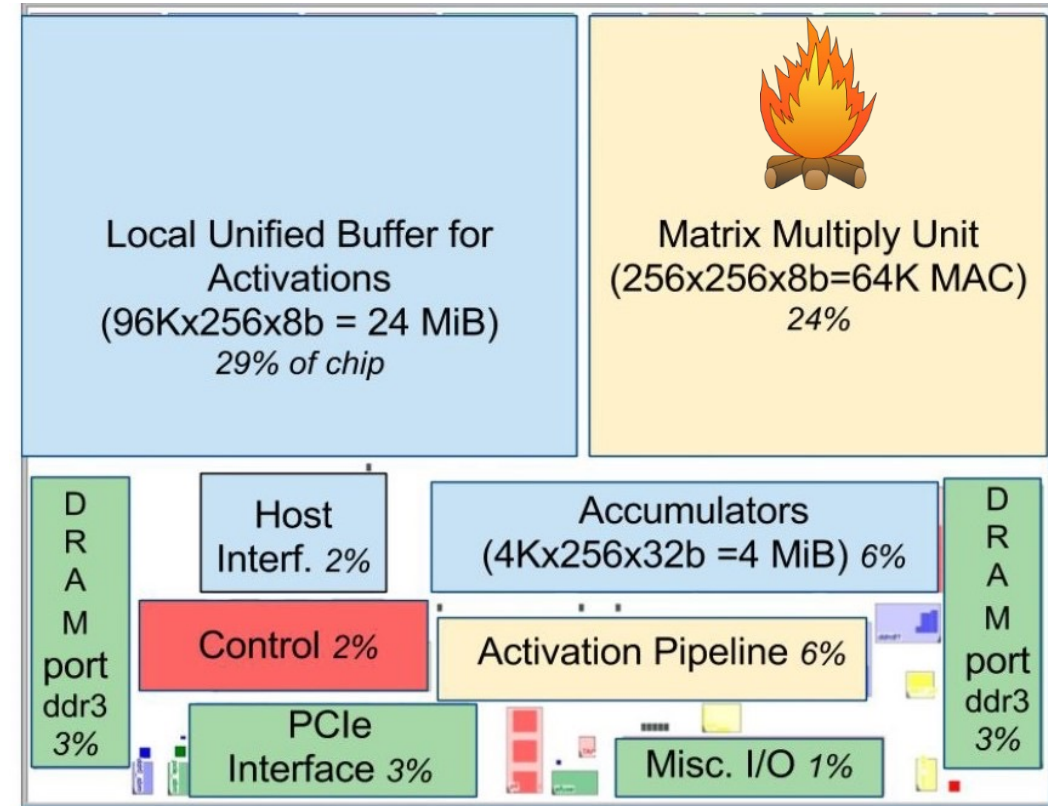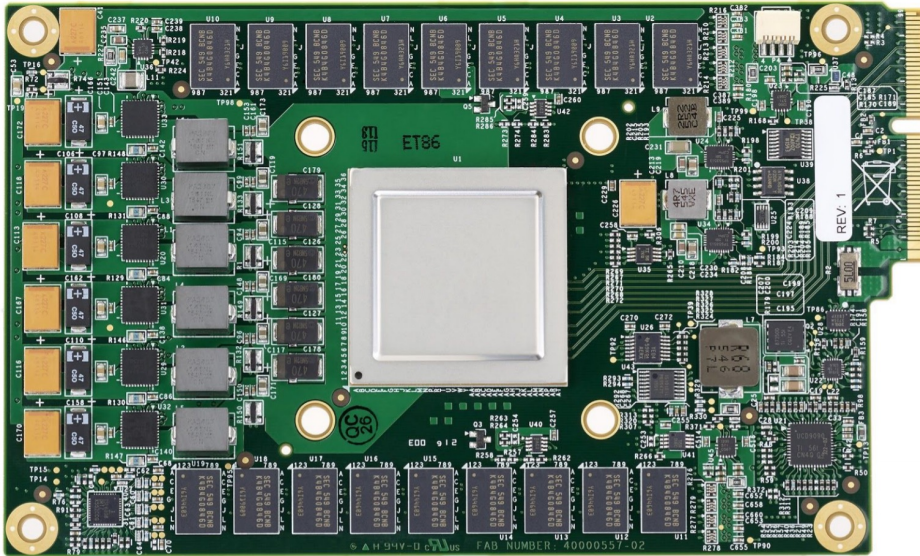# AI is reshaping the Future of Humankind

## *But At Which Cost?*



**Training** a single **AI** model emits carbon > **5x cars** in their lifetimes

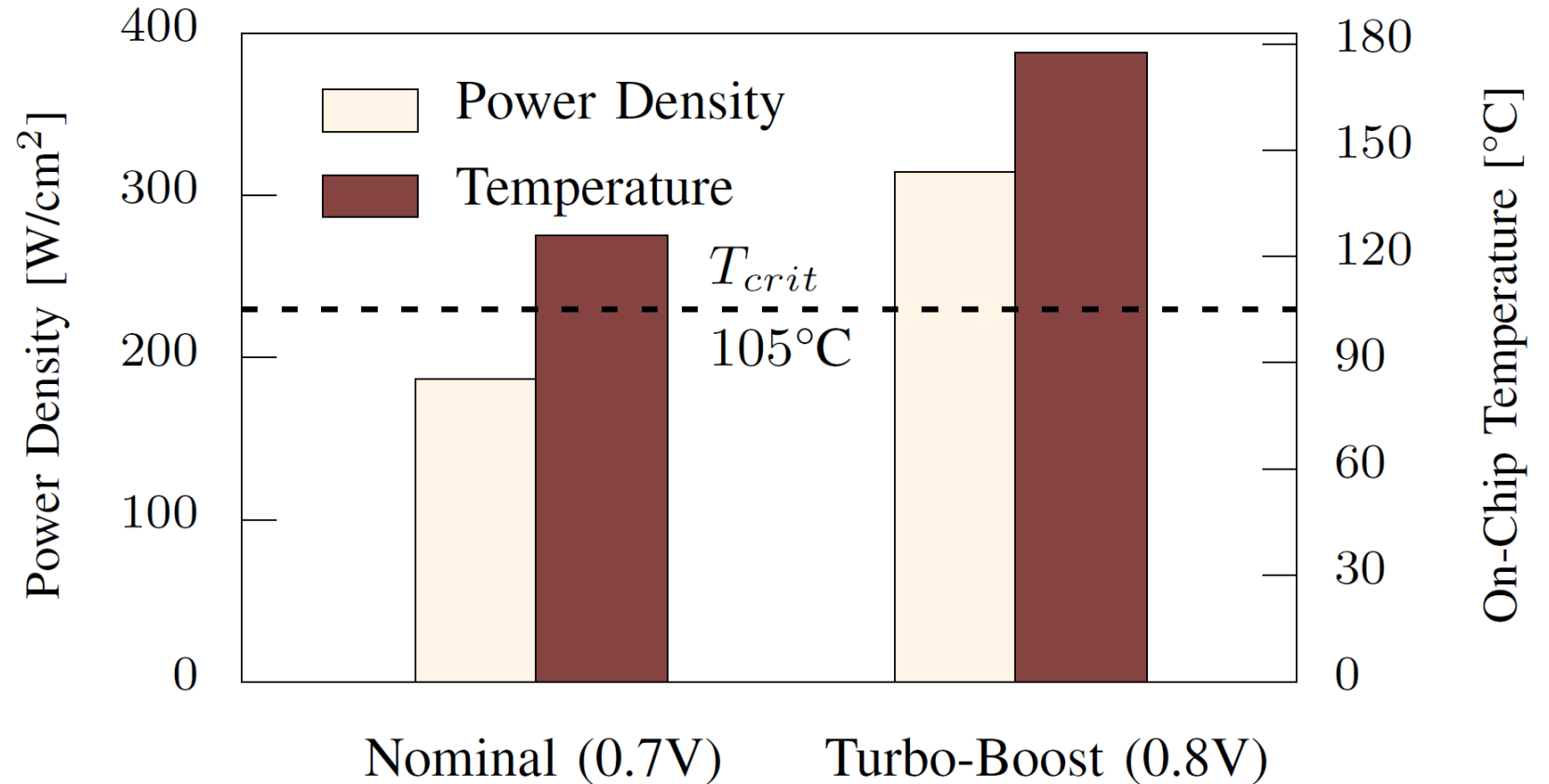src: Emma Strubell, et al. "Energy and Policy Considerations for Deep Learning in NLP" in 57th ACL, 2019.

$CO_2$

# Deep Learning is REALLY Power Hungary!



Google TPU [ISCA'17]

# Deep Learning is REALLY Power Hungry!

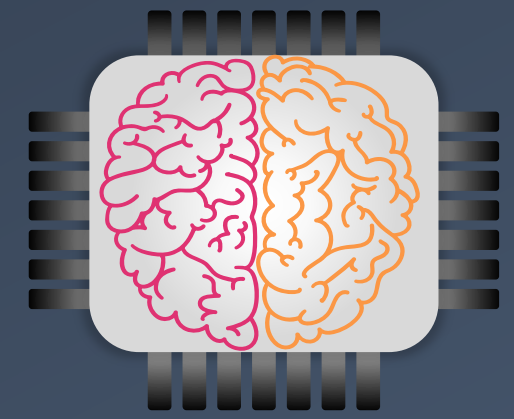**Why not alternative algorithm to Deep Learning?**
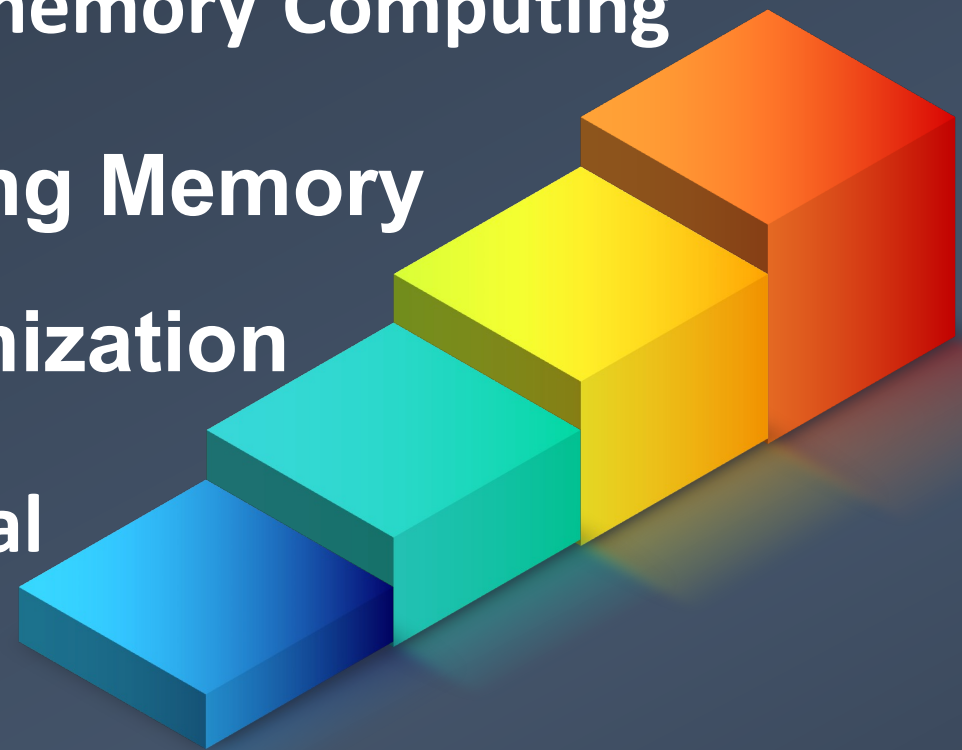


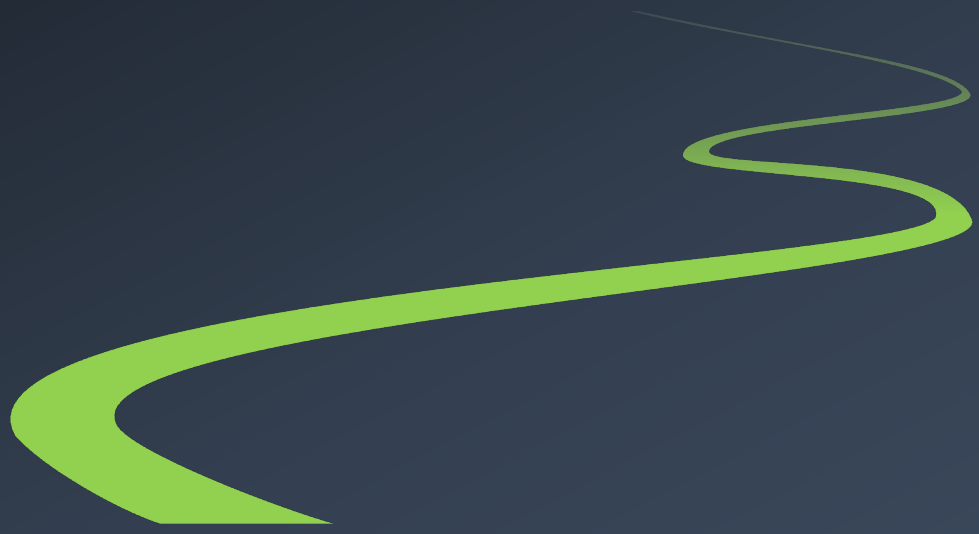Amrouch [TCAD'20]

**Hyperdimensional
in-memory Computing**

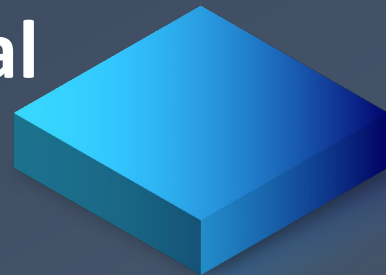**Emerging Memory**

**RISC-V Customization**

**Hyperdimensional
Computing**

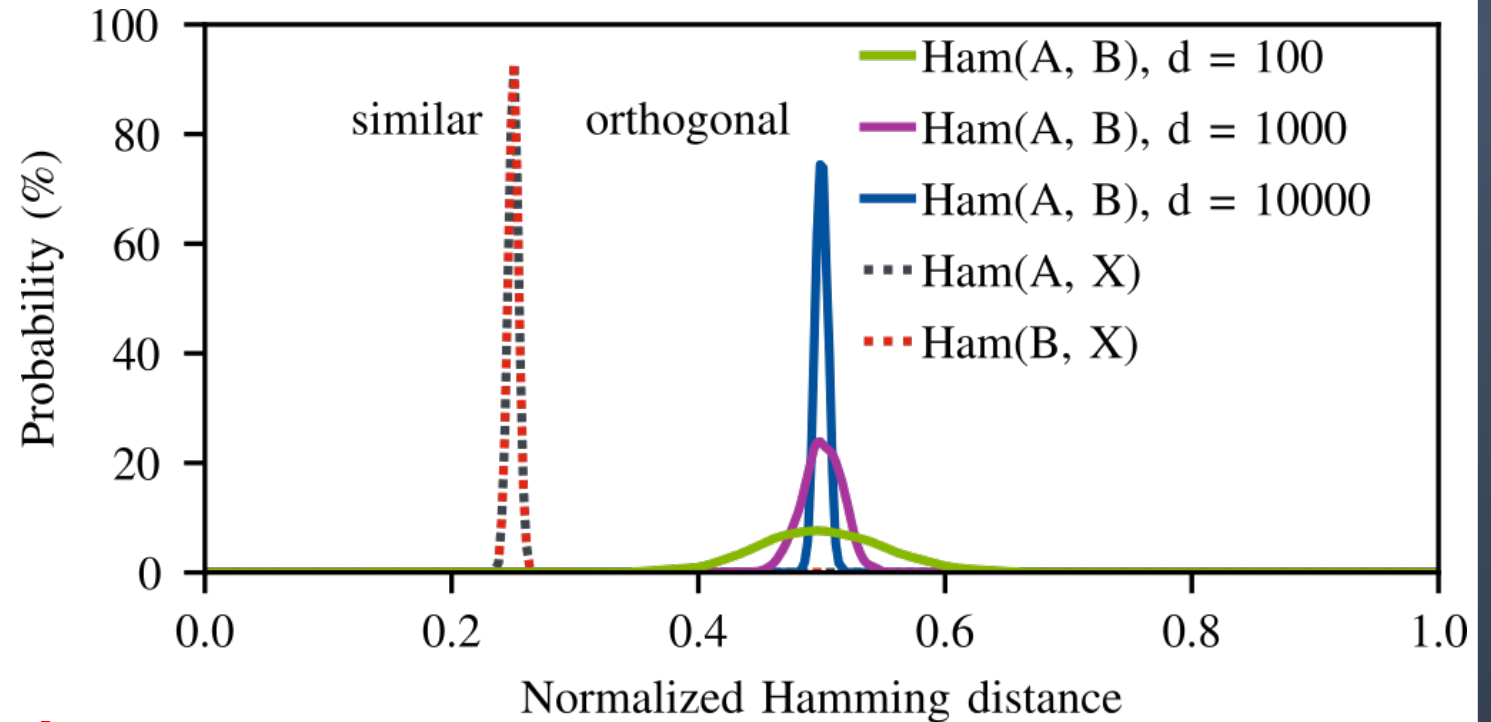**Brain-inspired
Computing
for Edge AI**

# Hyperdimensional Computing

# Brain-Inspired Hyperdimensional Computing

- Large vectors, e.g., 10000 elements

- Randomness is a feature not a bug

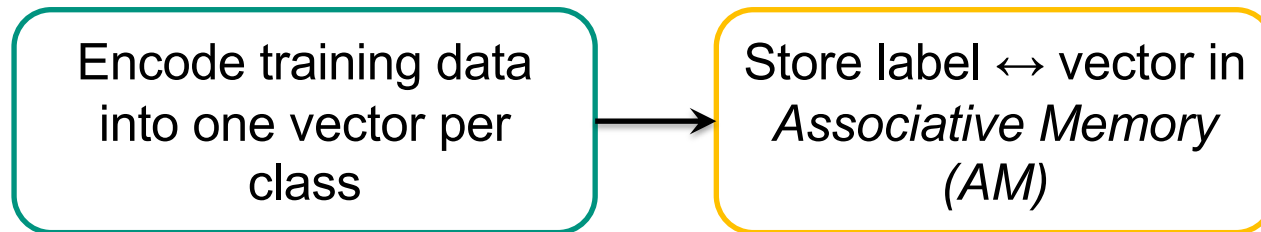- Simple Operations
  - Permutation
  - Binding
  - Bundeling



**Similarity is the Core Principle**
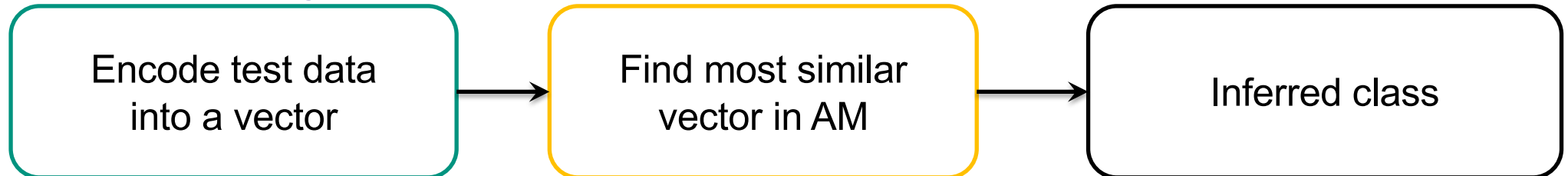
# Brain-Inspired Hyperdimensional Computing

## 1. Prepare: Encode real-world data into hyperspace

| Random vectors representing basic real-world values | → | Store in *Item Memory (IM)* | → | Combine item vectors to encode complex real-world data into a vector |
|---|---|---|---|---|

## 2. Learn: Train the model

| Encode training data into one vector per class | → | Store label ↔ vector in *Associative Memory (AM)* |
|---|---|---|

## 3. Inference: Recognize unknown data

| Encode test data into a vector | → | Find most similar vector in AM | → | Inferred class |
|---|---|---|---|---|

# Brain-Inspired Hyperdimensional Computing

**Example: Language classification**

(1) Assign a random vector: VERY large (10k bits)

```
a=[10110000010000110101]
b=[10100011011010000001]
         ⋮
!=[10101111000111100101]
```

(2) Encoding with N-Grams using two simple operations: **XOR**, **Rotate**

   **"Hi"  →  [H]  XOR  [Rotate(i)]**

# Brain-Inspired Hyperdimensional Computing

**Training Text:**     **To be, or not to be** ·········

[0010010000011111110001]

[0111010011010111111]

[1010010001011100010]

⋮

[1010010001011100010]

**Entire language is *just one Hyper Vector***

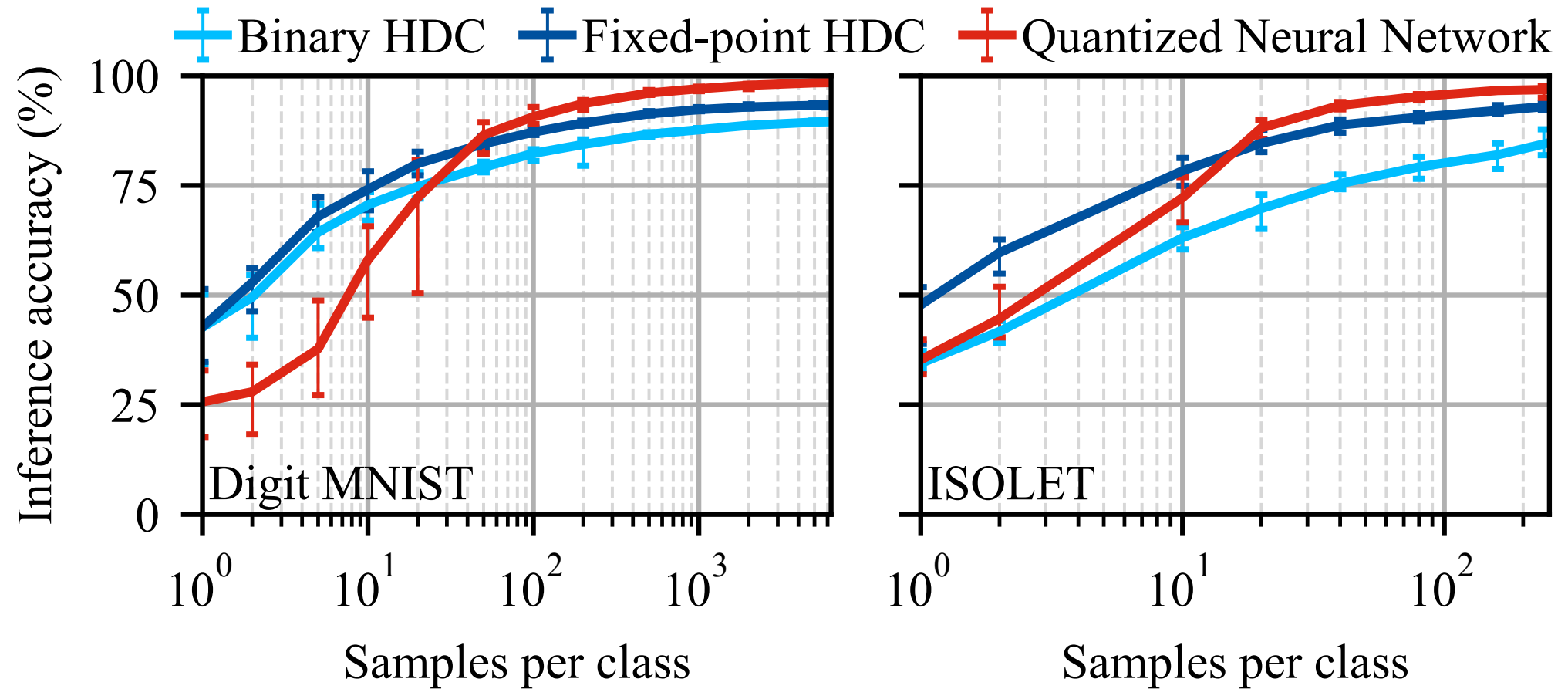Count 1's

[3,6,10,9,13,4,19,..70]

Majority gate ⟶ **[1010011000110011001]**

# Robustness against HW Errors and Noise
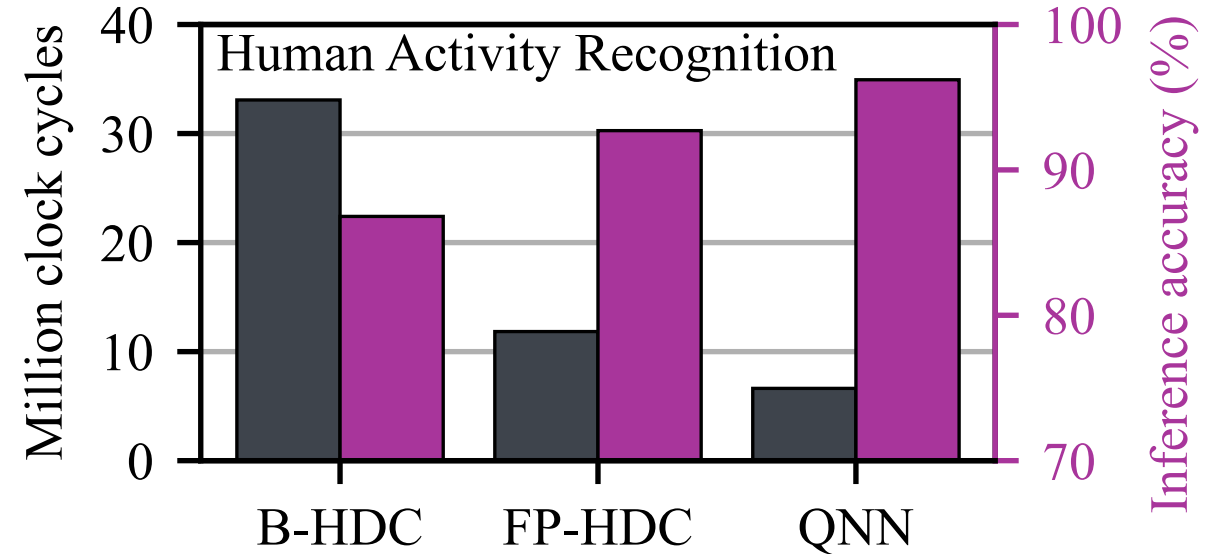


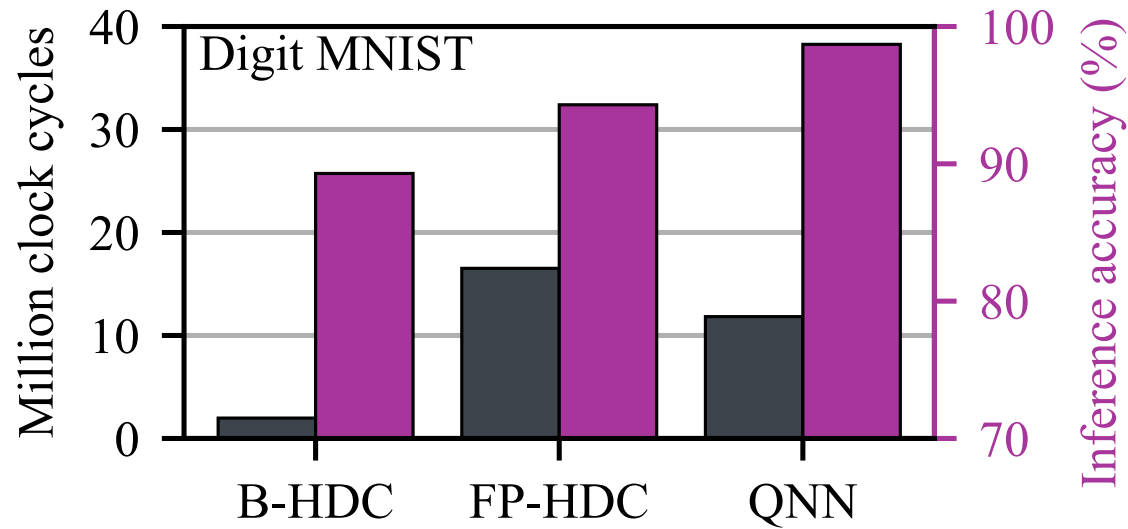Errors injected in the underlying HW operations

# HDC vs. QNNs: Learning from Little Data!



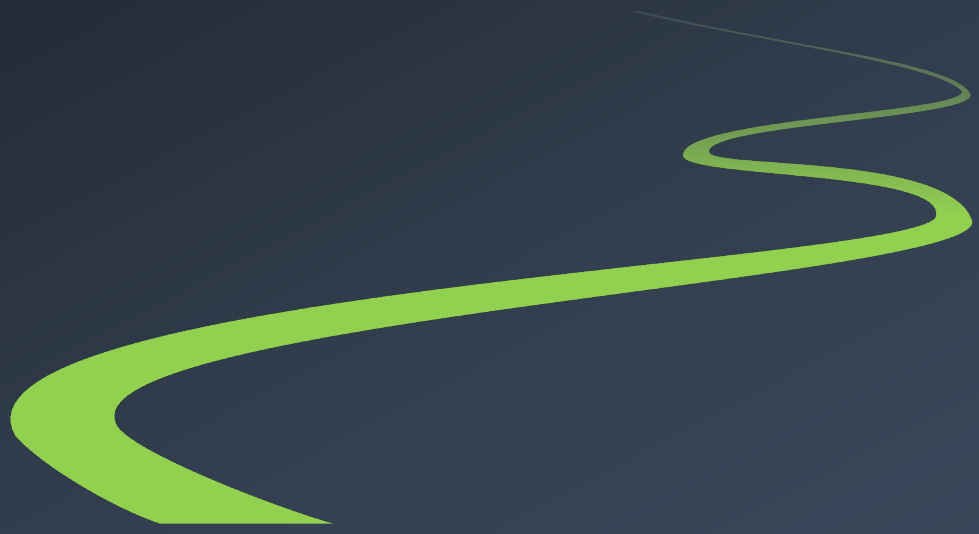HDC learns from little samples

Fix-point HDC : similar accuracy to QNN
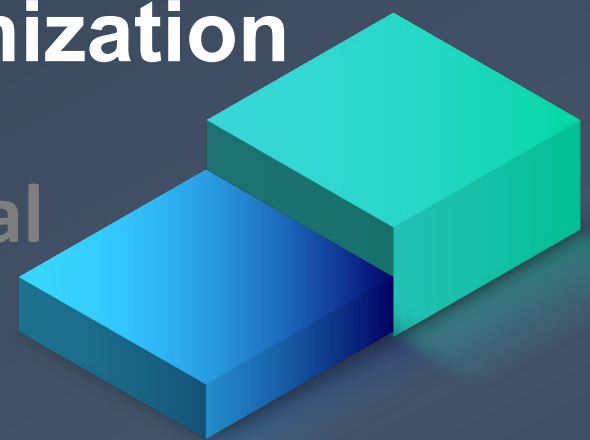
# HDC vs. QNN: Performance / Accuracy



Binary HDC has a superior speed in image classification

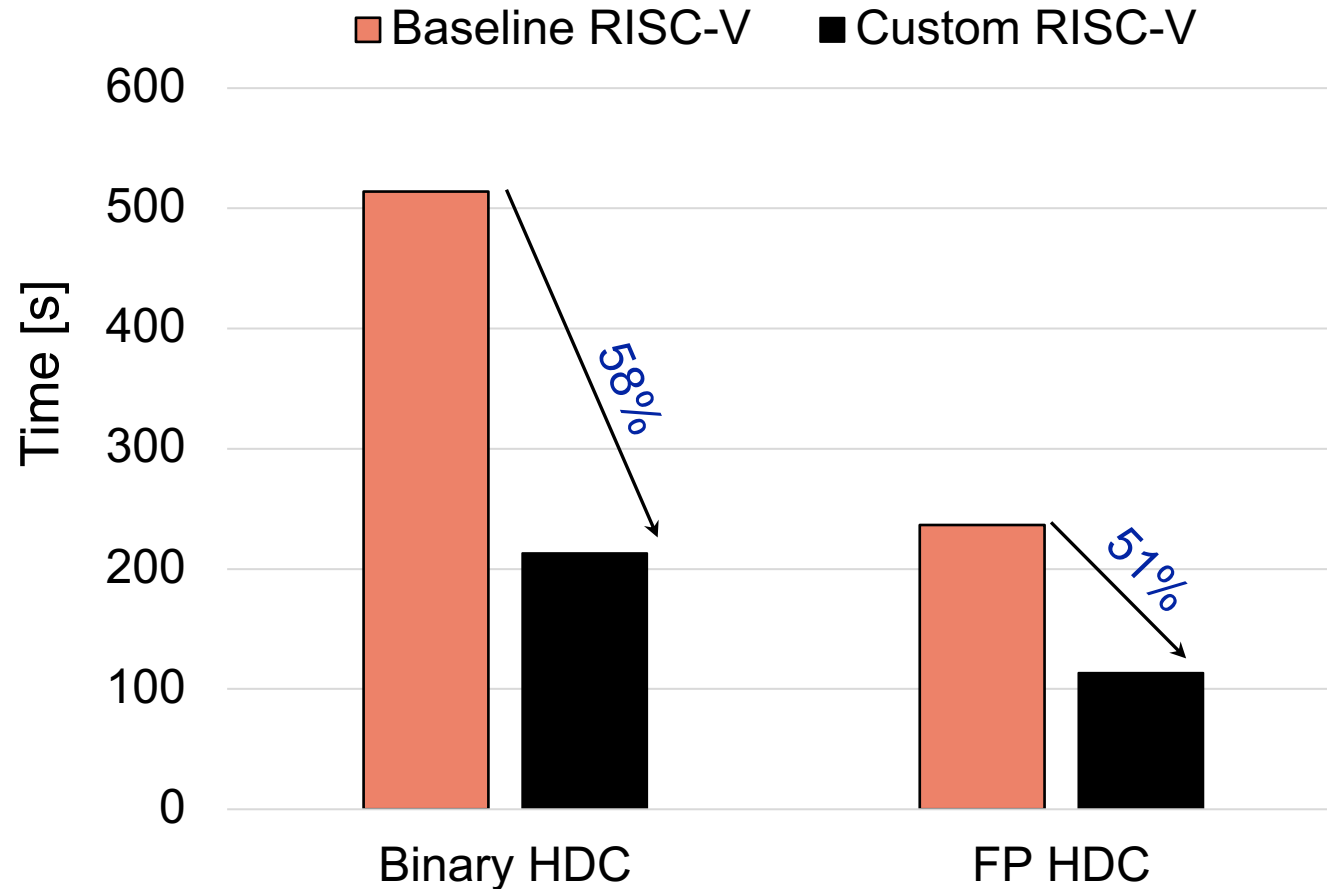Both QNN, HDC employ MACs, but QNN is faster than Fix-point HDC

# RISC-V Customization

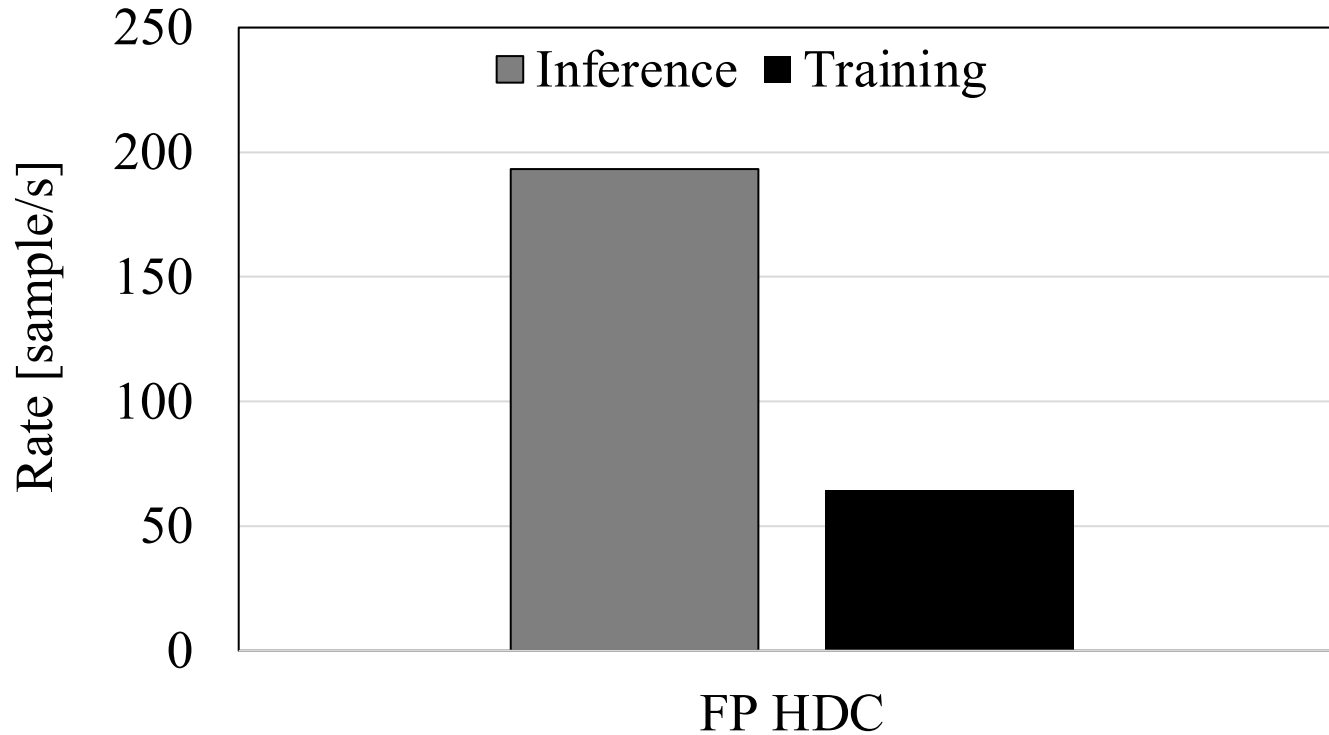## Hyperdimensional Computing

# RISC-V Customization for Edge AI: Training



Baseline RISC-V ■ Custom RISC-V

58%

51%

Binary HDC        FP HDC

SYNOPSYS® ASIP Designer

**Human activity dataset with ~7500 samples**

**Fully trained in less than 2 seconds!**

**Achieving a similar accuracy as QNN**

# RISC-V Customization for Edge AI: Training



Inference rate reaches
~200 samples per second

**SYNOPSYS**® **ASIP Designer**

# HDC: Inevitable Memory Bottleneck

**HDC relies on large vectors with > 1000 dimensions**

**→ Von Neumann architecture and memory bottleneck**

**In our analysis: Loading the vectors > 30% of cycles !**
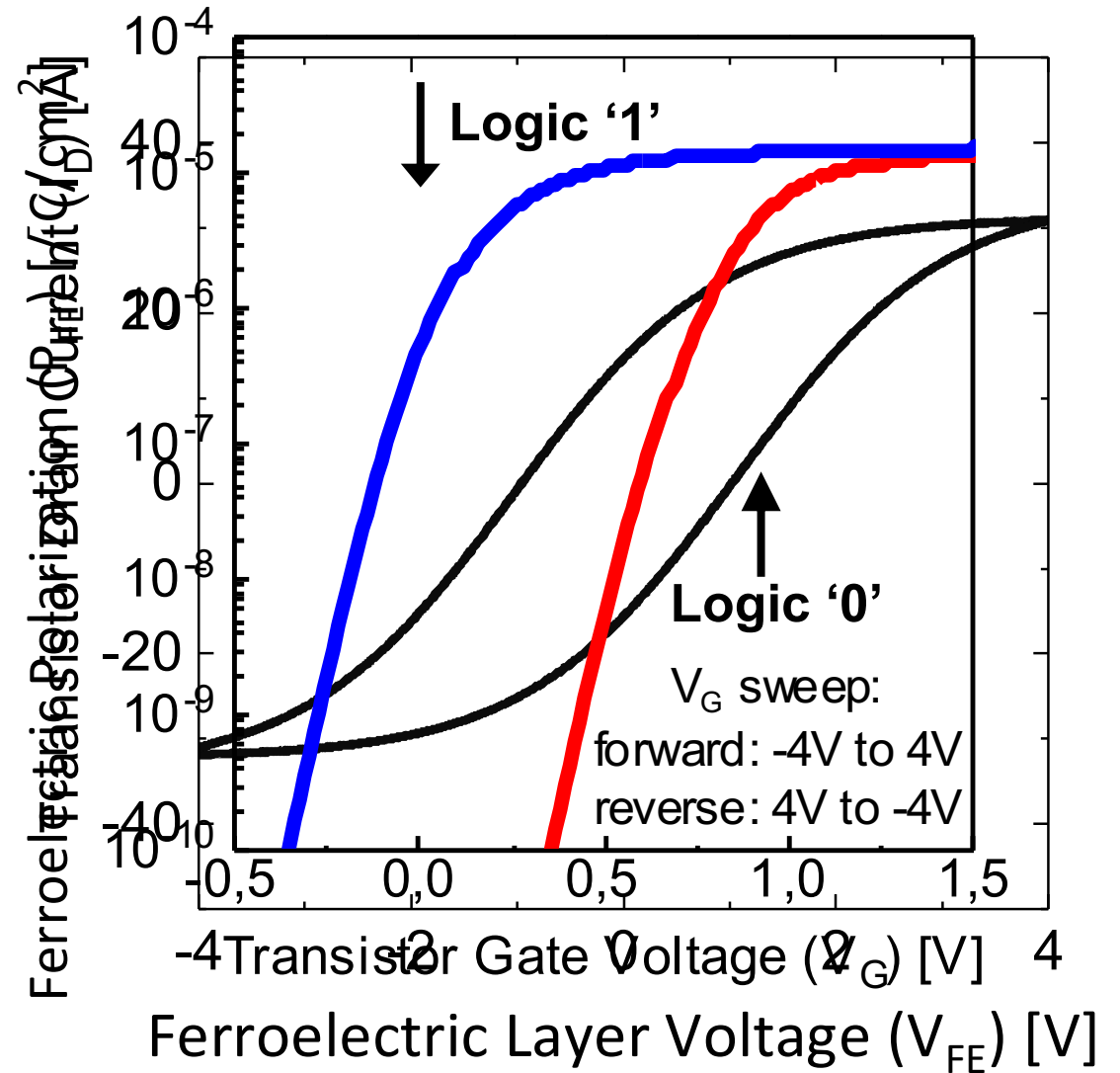
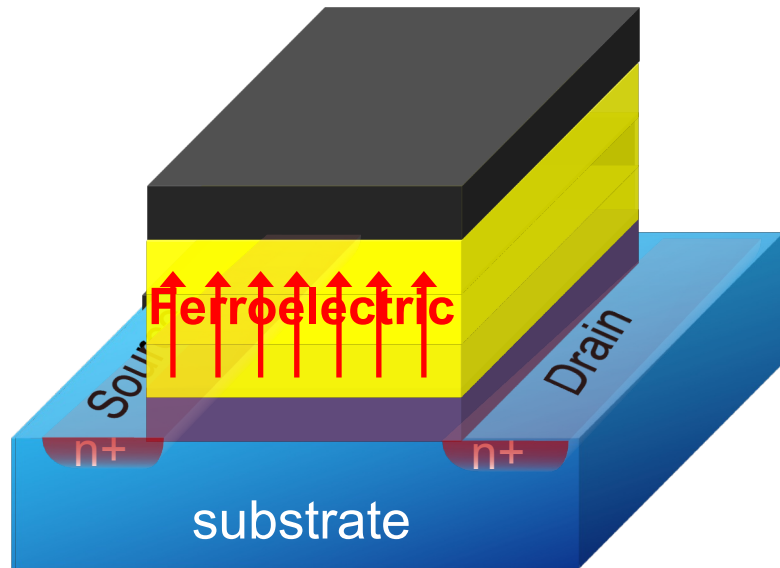**Hyperdimensional In-Memory Computing**

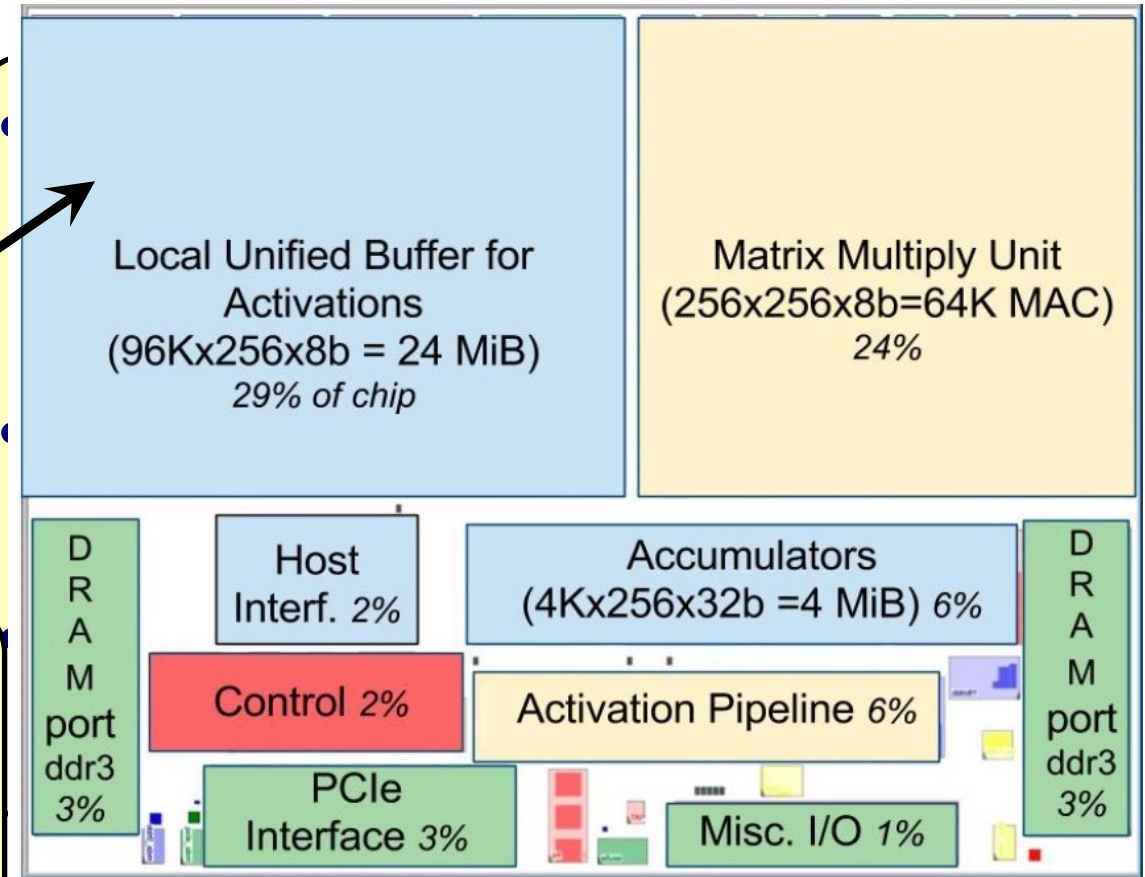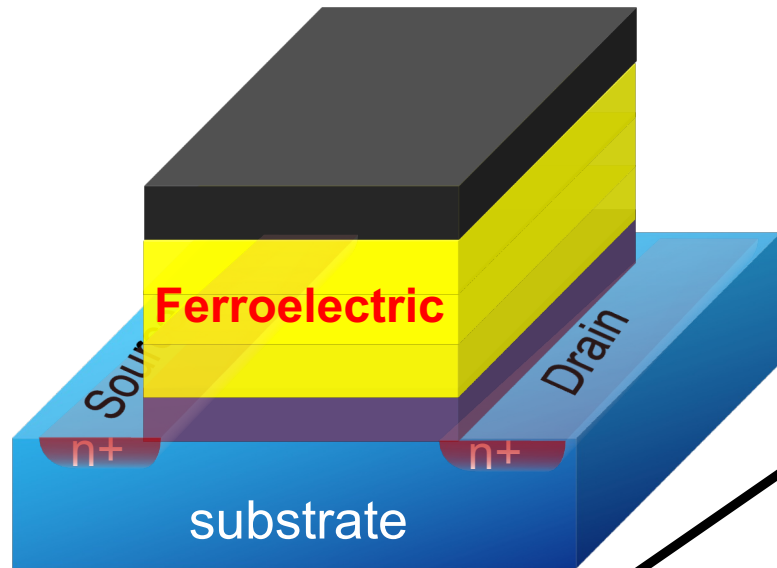**Emerging FeFET**

RISC-V Customization

Hyperdimensional Computing

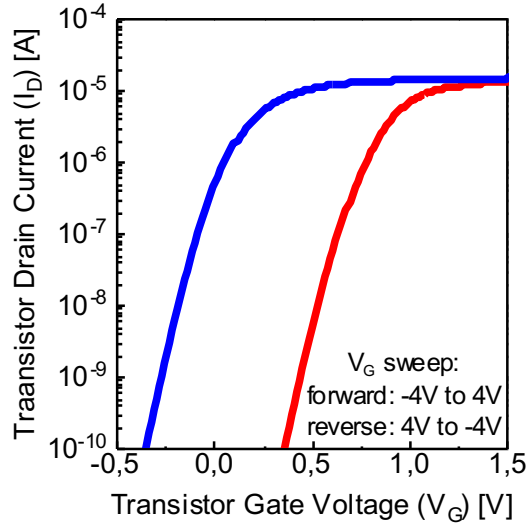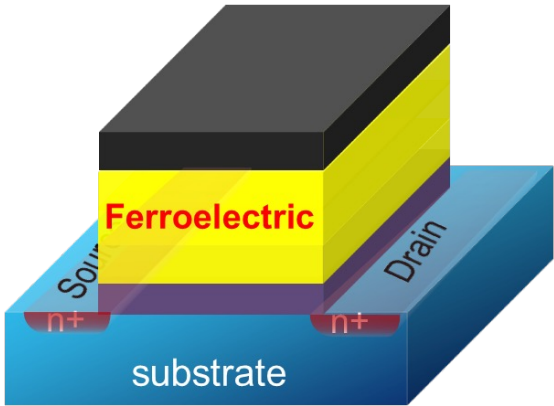**Brain-inspired Computing for Edge AI**

# From FET to FeFET

# FeFET: Emerging Memory

**Ferroelectric**

Source

Drain

n+    n+

substrate

**Replacing SRAMs with FeFETs**

→ **Large Power Saving**

→ **Higher Storage Capacity**

→ **Less DRAM Communications**

Local Unified Buffer for Activations
(96Kx256x8b = 24 MiB)
*29% of chip*

Matrix Multiply Unit
(256x256x8b=64K MAC)
*24%*

DRAM port ddr3 3%

Host Interf. 2%

Accumulators
(4Kx256x32b =4 MiB) 6%

Control 2%

Activation Pipeline 6%

PCIe Interface 3%

Misc. I/O 1%

DRAM port ddr3 3%

AI Chip: Google TPU [ISCA'17]

# In-Memory Computing using FeFETs



High current → '1'

Low current → '0'

V_G sweep:
forward: -4V to 4V
reverse: 4V to -4V

Vdd

SL — '0' — $\overline{SL}$

ON    OFF
OFF   OFF

Discharge
No Discharge

**Sensing the current**
→ *Match* or *Mismatch*

**Search '1'** →
**Search '0'** →

# In-Memory Computing using FeFETs



High current → '1'

Low current → '0'

ML

SL    $\overline{SL}$

Store '110'    '1'  '1'  '0'  '0'  '1'  '1'  '1'  '0'  '0'

Search '101'

match    mis    mis

**Hamming distance of 2**

H. Amrouch @ TUM Venture
Labs, E-mail: amrouch@tum.de

Chair of AI Processor Design,
Technical University of Munich

# In-Memory Hyperdimensional Computing

The row with the smallest hamming distance → *best match*



[10110000010000110101]

ML

[10110000010000110101]

ML

[10101111000111100101]

ML

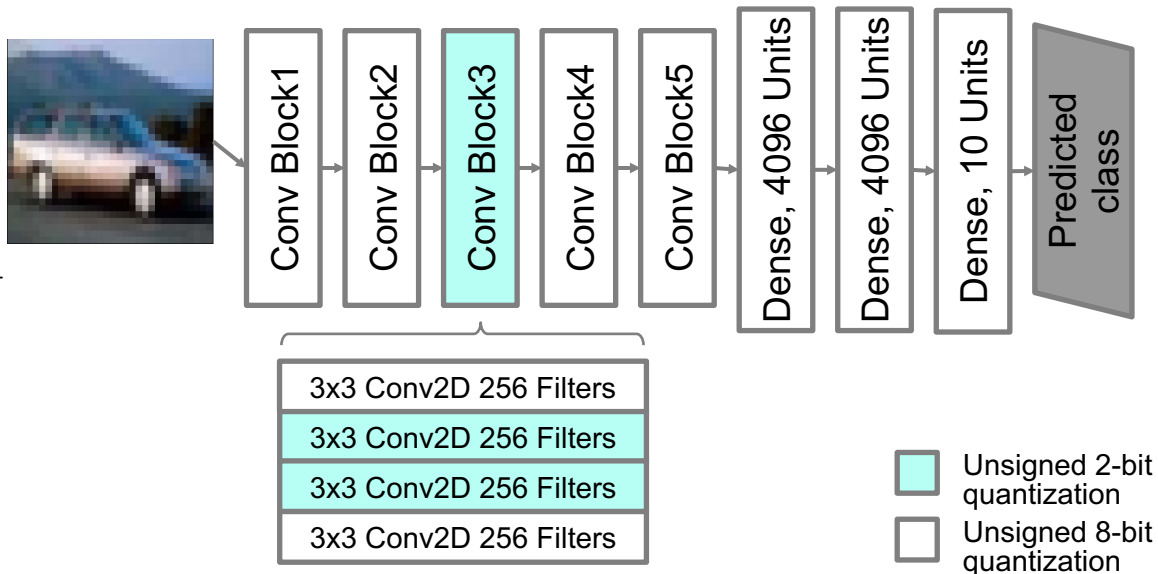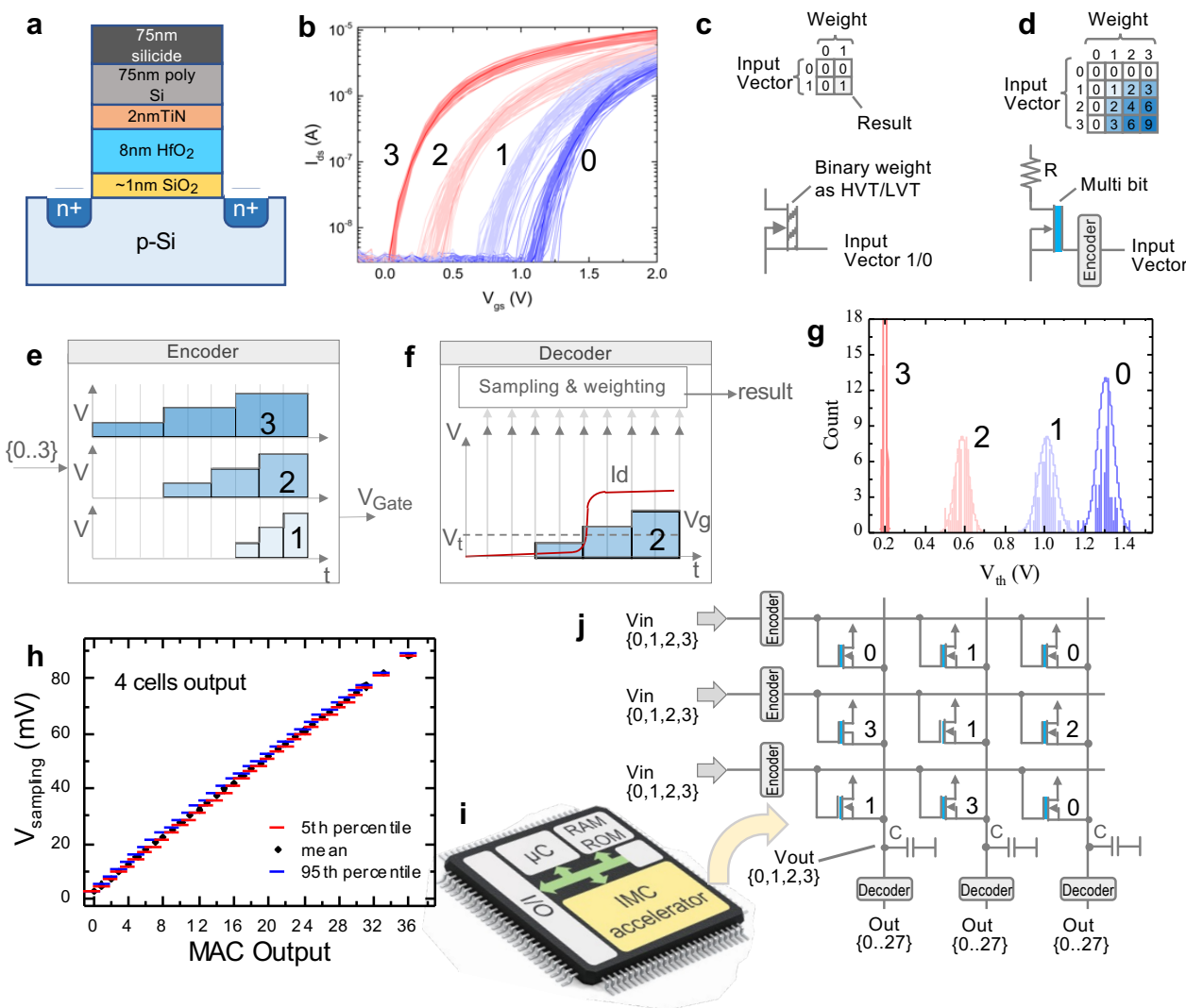10k cells

Search = [10010101100010001001001101]

S. Thomann, C. Li, C. Zhuo, O. Prakash, X. Yin, X. S. Hu, and H. Amrouch, "On the Reliability of In-memory Computing: Impact of Temperature on Ferroelectric TCAM," **IEEE VLSI Test Symposium (VTS'21), 2021** (Best Paper Nomination)
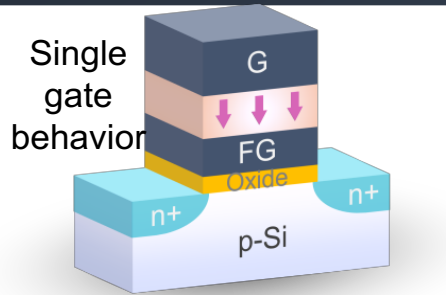
# Very Efficient MAC using FeFET Crossbar
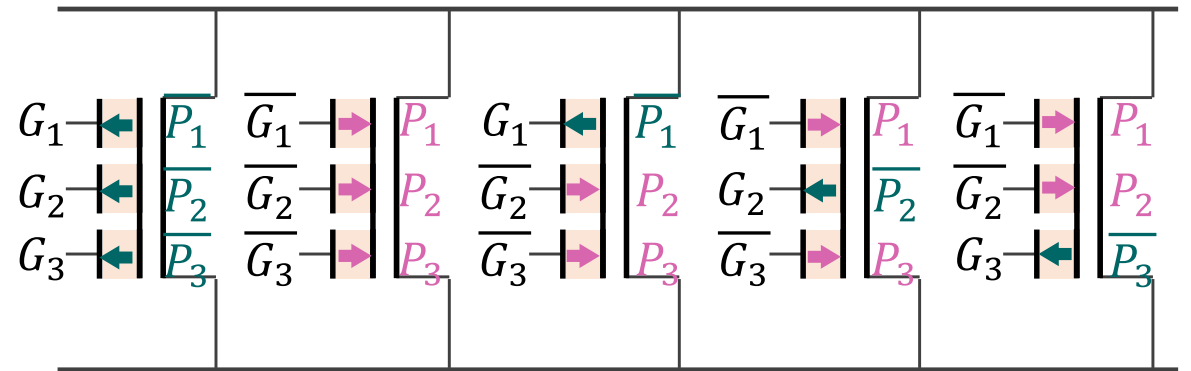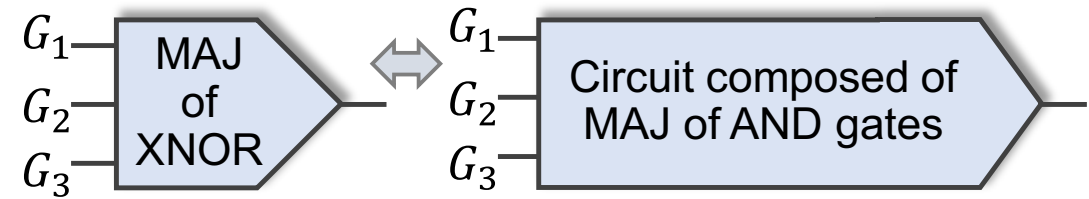
# From In-Memory → In-Transistor Computing

Single gate behavior



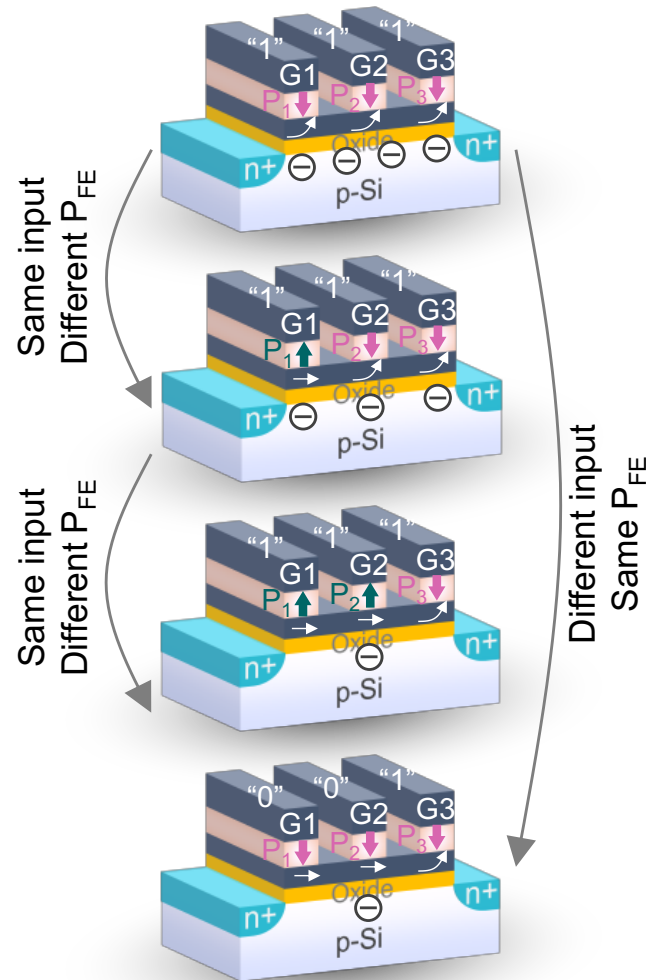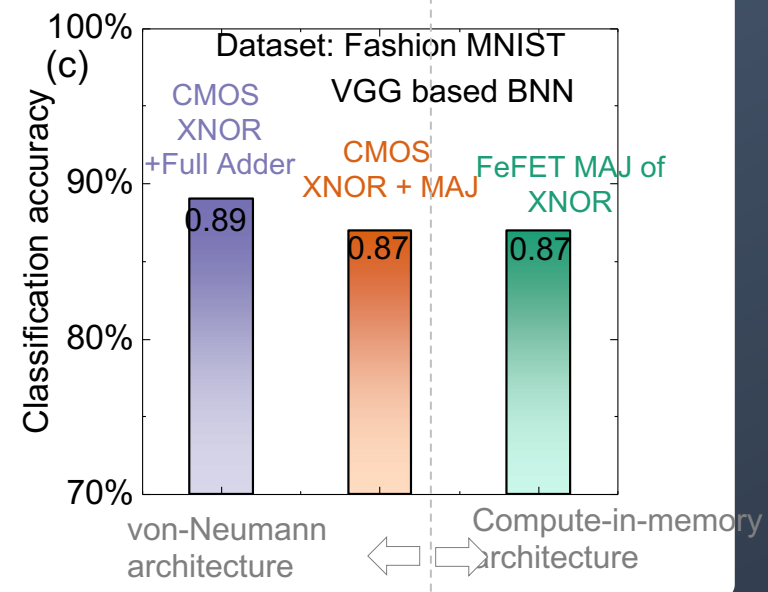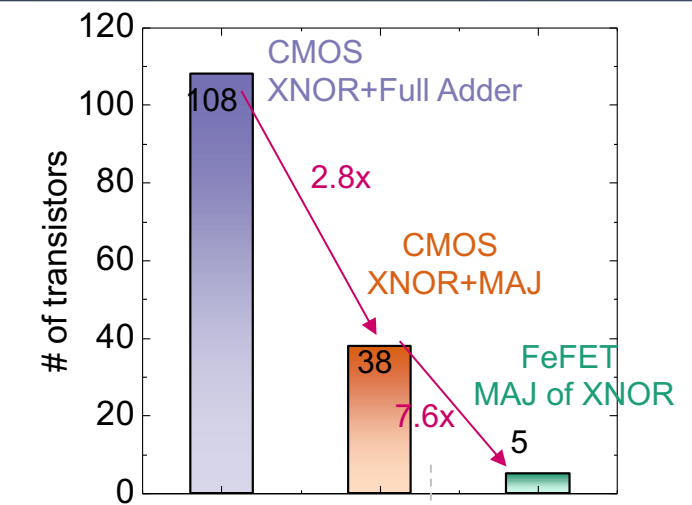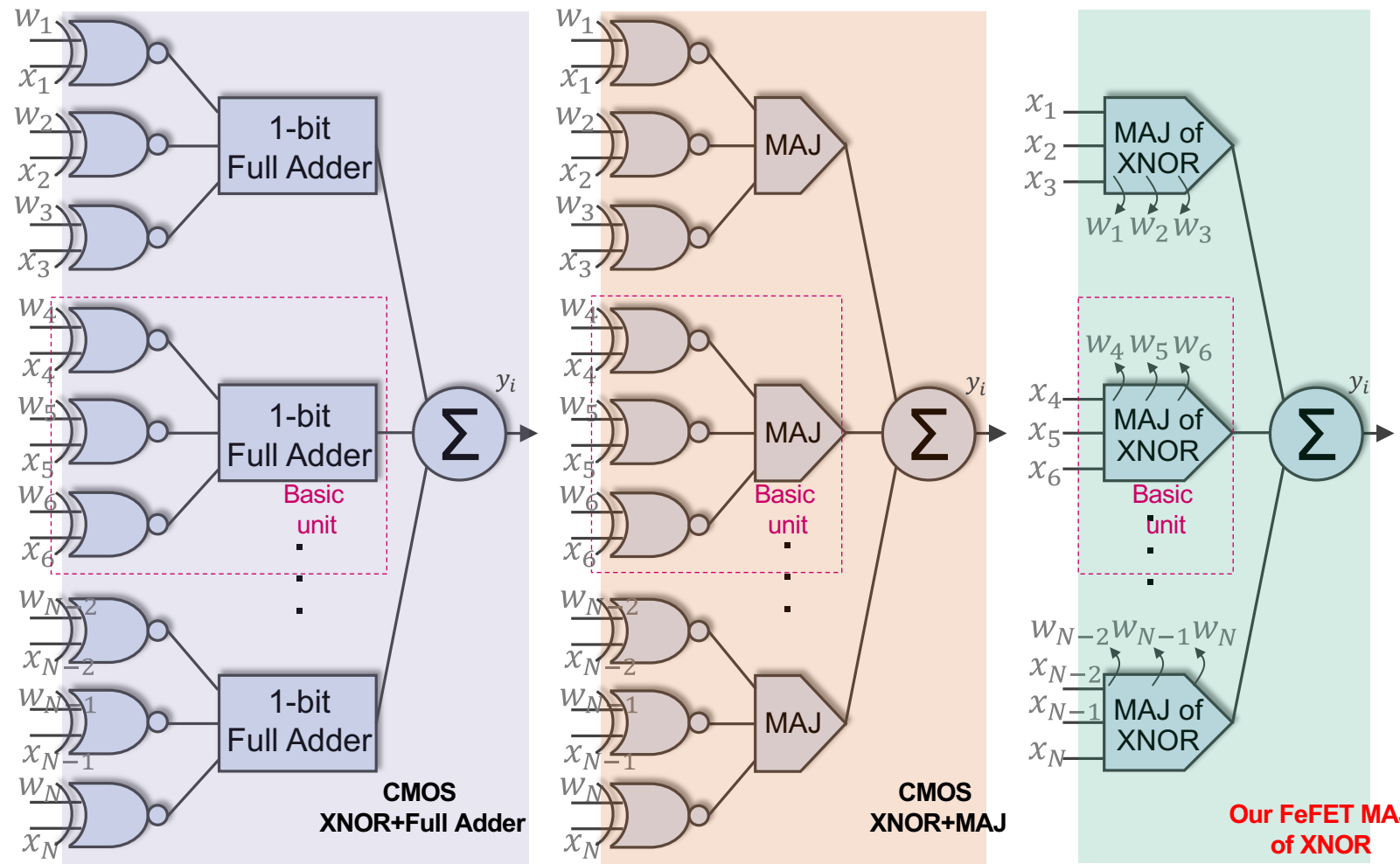Each gate $i$ performs AND operation: $G_i \wedge P_i$



$G_i$="0"    $G_i$="1"

Contribution of gate $i$ to $V_{FG}$

| $G_i$ \ $P_i$ | "1" ↓ | "0" ↑ |
|---|---|---|
| "1" (e.g., 1 V) | High | Low |
| "0" (e.g., -1 V) | Low | Low |

Same input Different $P_{FE}$

Same input Different $P_{FE}$

Different input Same $P_{FE}$



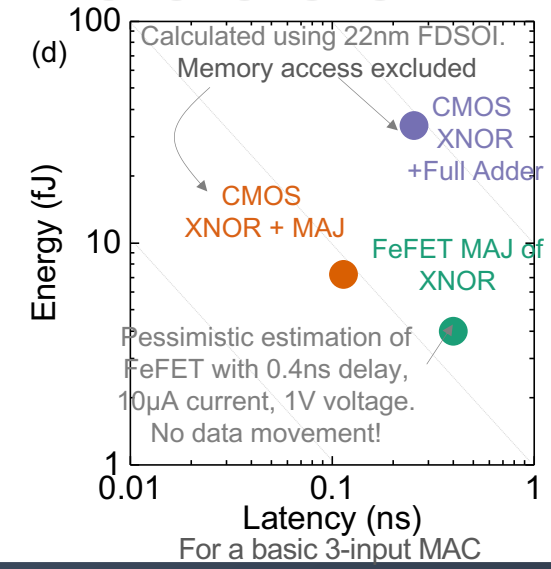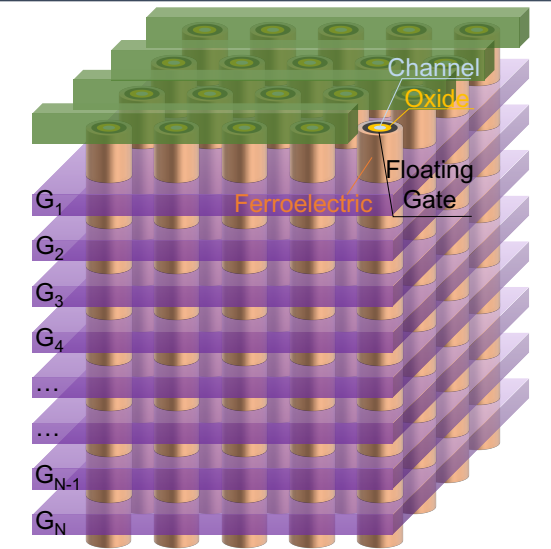$G_1$ — MAJ of XNOR — ⟺ — $G_1$, $G_2$, $G_3$ — Circuit composed of MAJ of AND gates



S. Thomann / H. Amrouch, "Compact ferroelectric programmable majority gate for compute-in-memory applications," in 68th Annual IEEE International Electron Devices Meeting (IEDM), 2022
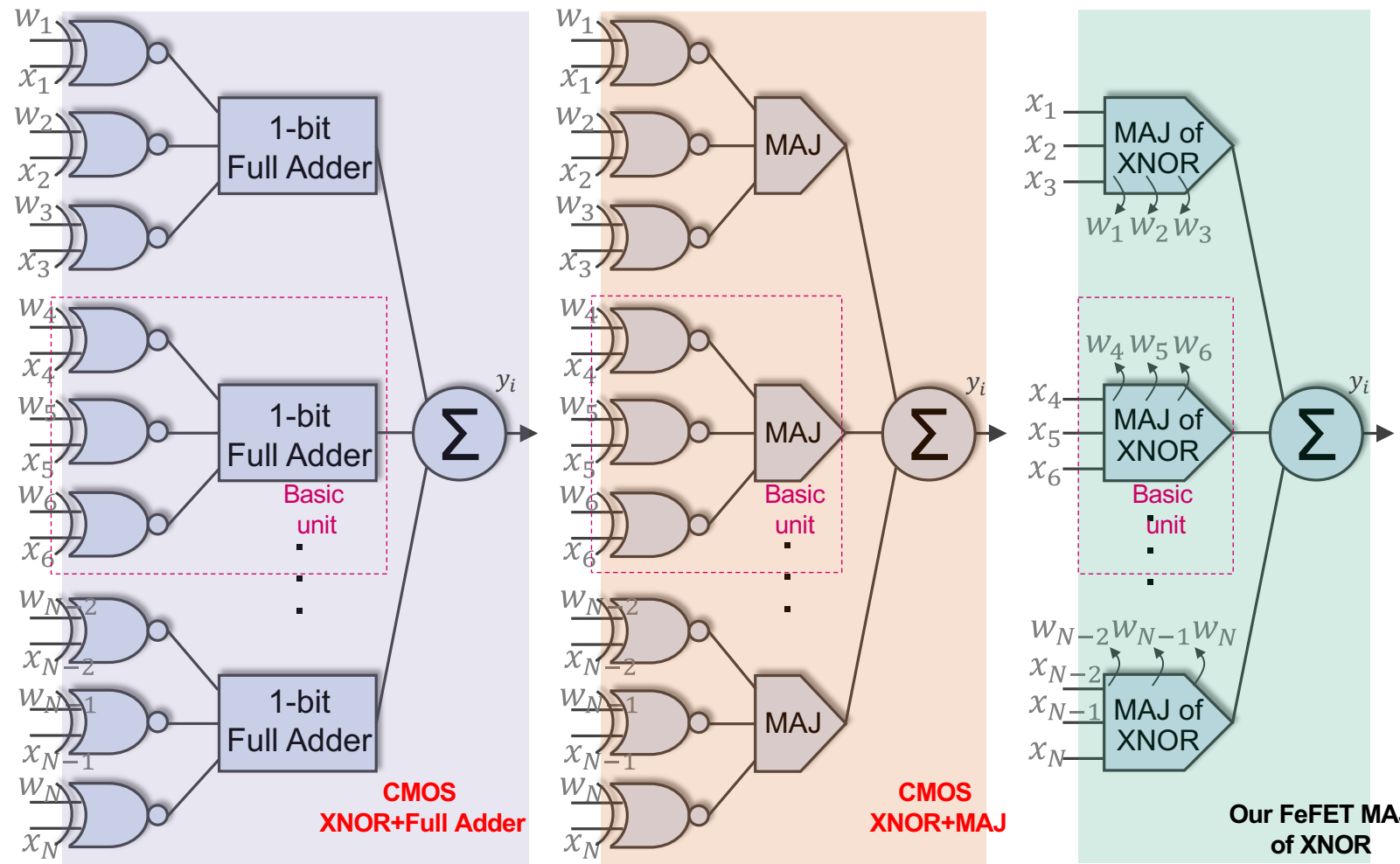
S. Thomann / H. Amrouch, "Compact ferroelectric programmable majority gate for compute-in-memory applications," **in 68th Annual IEEE International Electron Devices Meeting (IEDM), 2022**
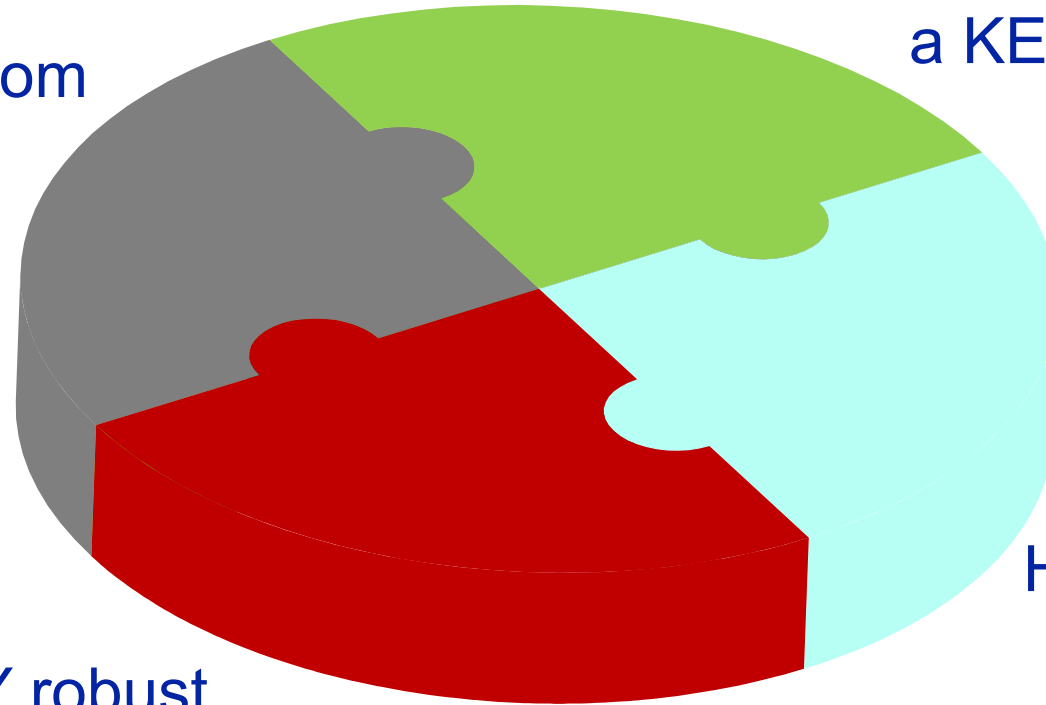
# From In-Memory → In-Transistor Computing



S. Thomann / H. Amrouch, "Compact ferroelectric programmable majority gate for compute-in-memory applications," in 68th Annual IEEE International Electron Devices Meeting (IEDM), 2022

# HDC Computing … Hope or Hype?



HDC can learn from little data

In-Memory Computing is also a KEY for very efficient HDC

HDC is VERY robust against errors

HDC enables training on the edge BUT RISC-V customization is the KEY

# Without them, nothing would be possible
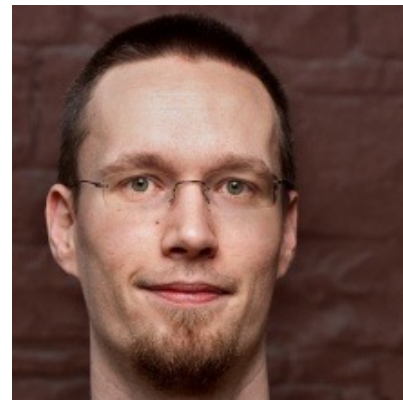
## Device / Circuit



Shubham Kumar (PhD)



Swetaki Chatterjee (PhD)

## Device / Circuit



Shivendra Parihar (PhD)



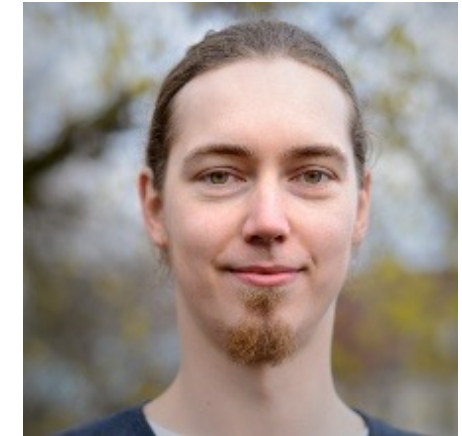Dr. Victor van Santen

## Digital Design



Shubham Kumar (PhD)



Simon Thomann (PhD)

## Deep Learning



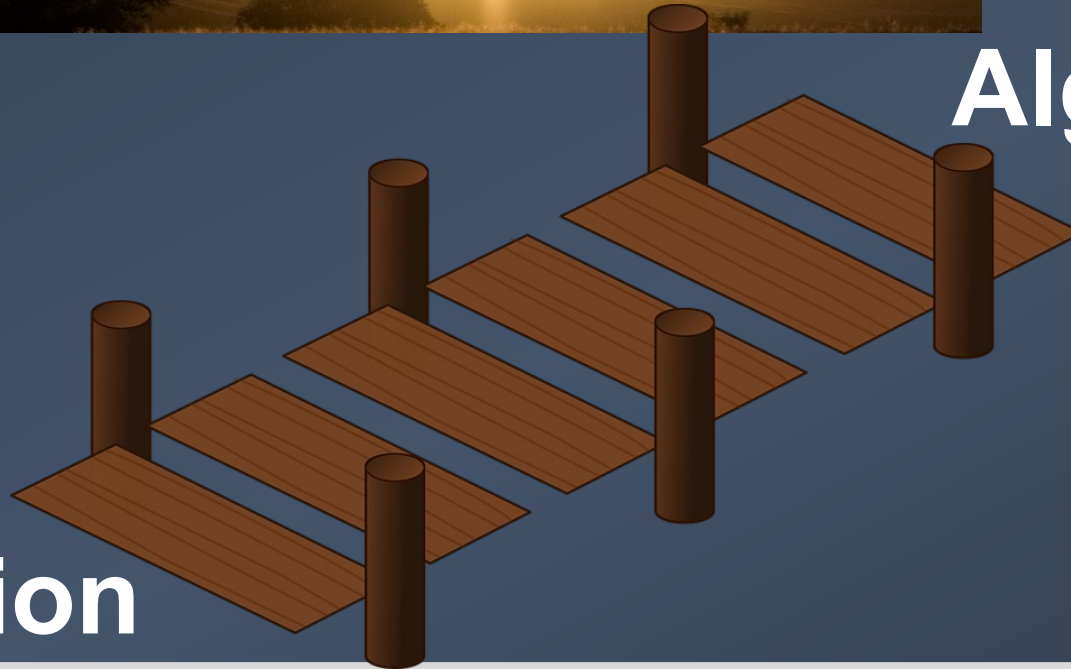Paul Genssler (PhD)



Rodion Novkin (PhD)

# Acknowledgement

# On the Brink of a new Era in Edge AI



**Novel AI Algorithms**

**RISC-V Customization**