# TOWARDS NEXT-GENERATION COMPUTING



### IN-MEMORY COMPUTING TECHNOLOGIES FOR DEEP NEURAL NETWORK ACCELERATION D.VERKEST



#### **SMART SERVICES** POWERED BY MACHINE LEARNING IN DATA CENTER



IMEC

#### SMART SERVICES POWERED BY MACHINE LEARNING IN DATA CENTER



IMEC

#### TODAY'S DOMINANT SOLUTION TARGET HIGH-PERFORMANCE MARKET

- Deep Neural Networks
  - CNN, LSTM
  - Vision, voice assistant
- General purpose platforms support
  - Training and inference
  - Flexibility in still evolving field
- Too power hungry for edge devices

5



#### TODAY'S DOMINANT ML PLATFORMS NOT FIT FOR SMART THINGS

- Intelligence at extreme edge
  - Support high performance DNN inference
  - Always-on
  - On a tight power budget
- → I to I0 PetaOP/s/W



10,000 TOPISIN



PUBLIC











PUBLIC

IMEC















- Vin Rw Vin 2 Rw Vin\_3 Rw Vin\_4 Rw  $I \sim \sum_{i=0}^{N} x_i w_{ii}$ PUBLIC IMEC
- Use memory array for massive parallel analog implementation of multiply-accumulate operations in DNN layer
- Memory array stores weights <u>and</u> implements the neural network layer matrix-vector multiplication in analog fashion
  - ightarrow analog compute-in-memory
  - ightarrow computational memory
  - ightarrow neuromorphic computing

#### ANALOG CIM ACCELERATORS FOR ML PROMISING SOLUTION FOR ENERGY EFFICIENCY

- I0,000 TOP/s/W
  - Memory cell
    0.25µm x 0.25µm
  - I6M cells/mm<sup>2</sup> at I00MHz
  - I.6 PMAC/s/mm<sup>2</sup> for ~0.2fJ/MAC

65th

International

Electron Devices

Meetina



IMEC

10mm<sup>2</sup> 20,000<sup>T0P5/W</sup>

Tue Dec 10, 2019, Focus session: Emerging AI hardware, Paper 22.2 by S. Cosemans, et al.

#### HOW TO MAKE IT WORK? SYSTEM-TECHNOLOGY CO-OPTIMIZATION





- Deep neural networks
  - Perform inference with lower precision
  - Can tolerate some weight variation
- Architecture
  - Use multiple memory cells per weight for more precision
- Circuit and device
  - Tune devices for neural network operation

#### ANALOG COMPUTE-IN-MEMORY ACCELERATORS FOR ML PROOF-OF-CONCEPT MATRIX-VECTOR MULTIPLICATION (MVM)

10,000 TOPS/W 1,000 TOPS/W SRAM-based Log speed (GOP/s) compute cells for 1000000 MVM of CNN and ЫТ LSTM 1,300 TOP/s/W 0 8V 100000 ANIA 2.700 TOP/s/W 0.6V FDX22 10000 4mm2 incl. 0-100-0 500k cells • mi • · (()) Ternary weights 1000 7-bit activations D. Bankman ISSCC2018 6-bit ADC (MAC) 100 Record efficiencies up Based on https://nicsefc.ee.tsinghua.edu.cn/projects/neural-network-accelerator/ to 2,700 TOP/s/W 10 0.1 10 0.001 0.01 100 1000 Log (Power) (W)

# embracing a better life

## THANK YOU FOR YOUR ATTENTION

