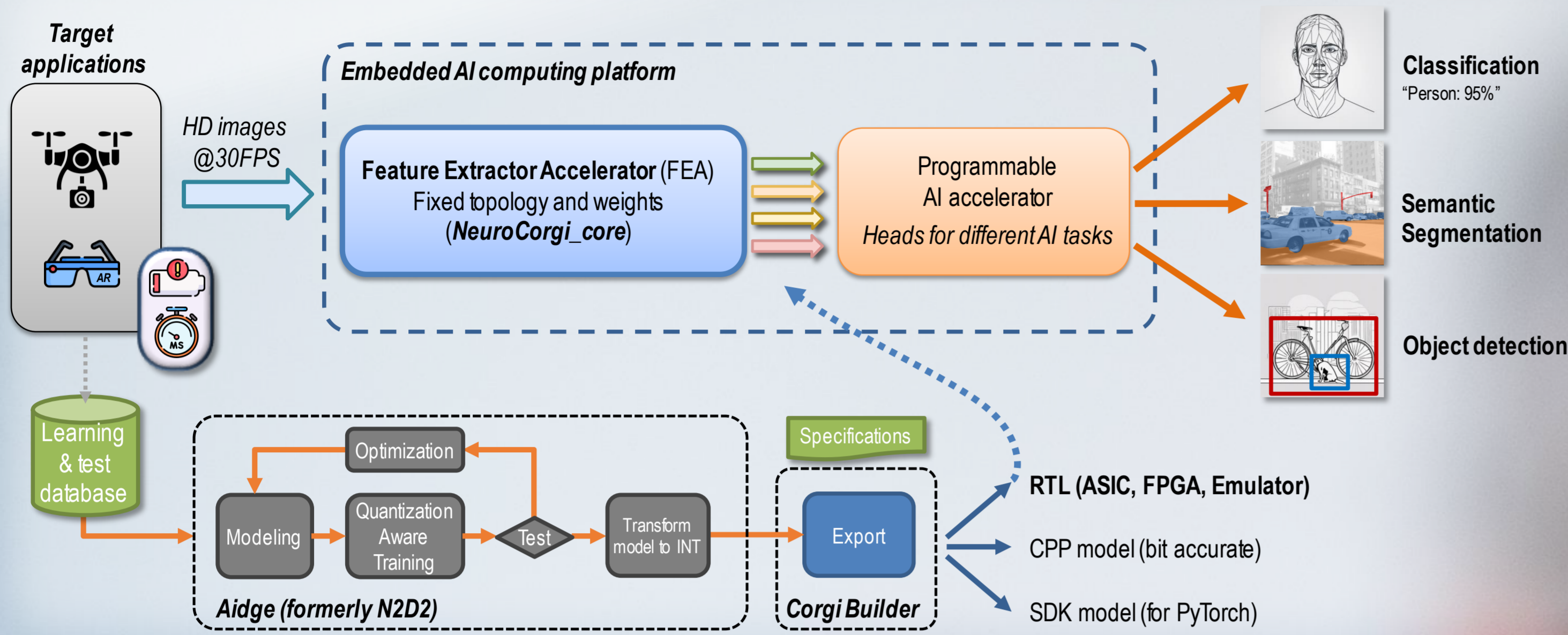


# A 772μJ/frame ImageNet Feature Extractor Accelerator on HD Images at 30FPS

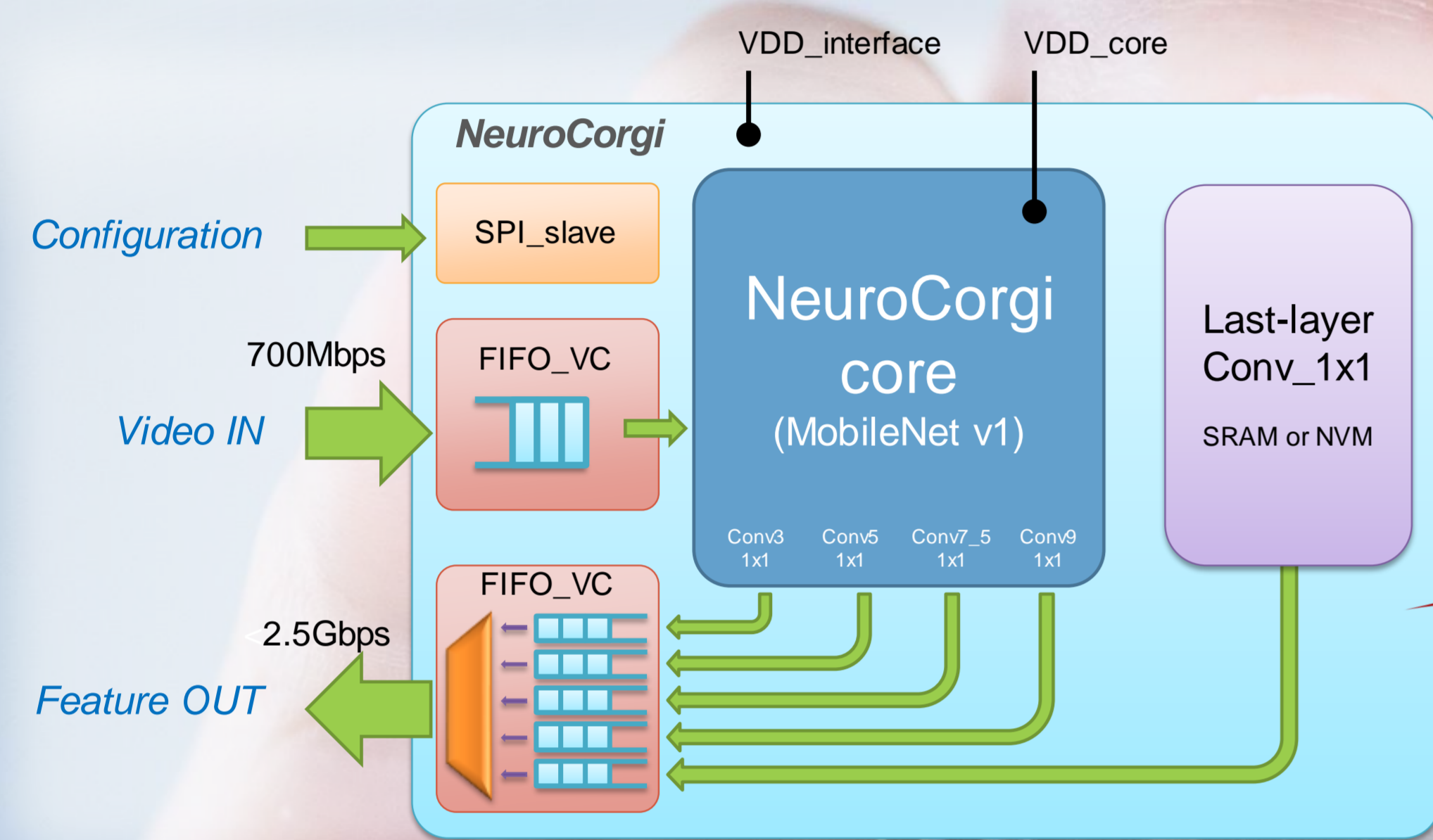
Ivan Miro-Panades<sup>1</sup>, Vincent Lorrain<sup>2</sup>, Lilian Billod<sup>1</sup>, Inna Kucher<sup>2</sup>, Vincent Templier<sup>2</sup>, Sylvain Choisnet<sup>1</sup>, Nermine Ali<sup>2</sup>, Baptiste Rossignaux<sup>2</sup>, Olivier Bichler<sup>2</sup>, Alexandre Valentian<sup>1</sup>  
<sup>1</sup>CEA, List, Grenoble, France  
<sup>2</sup>CEA, List, Palaiseau, France

## Overview

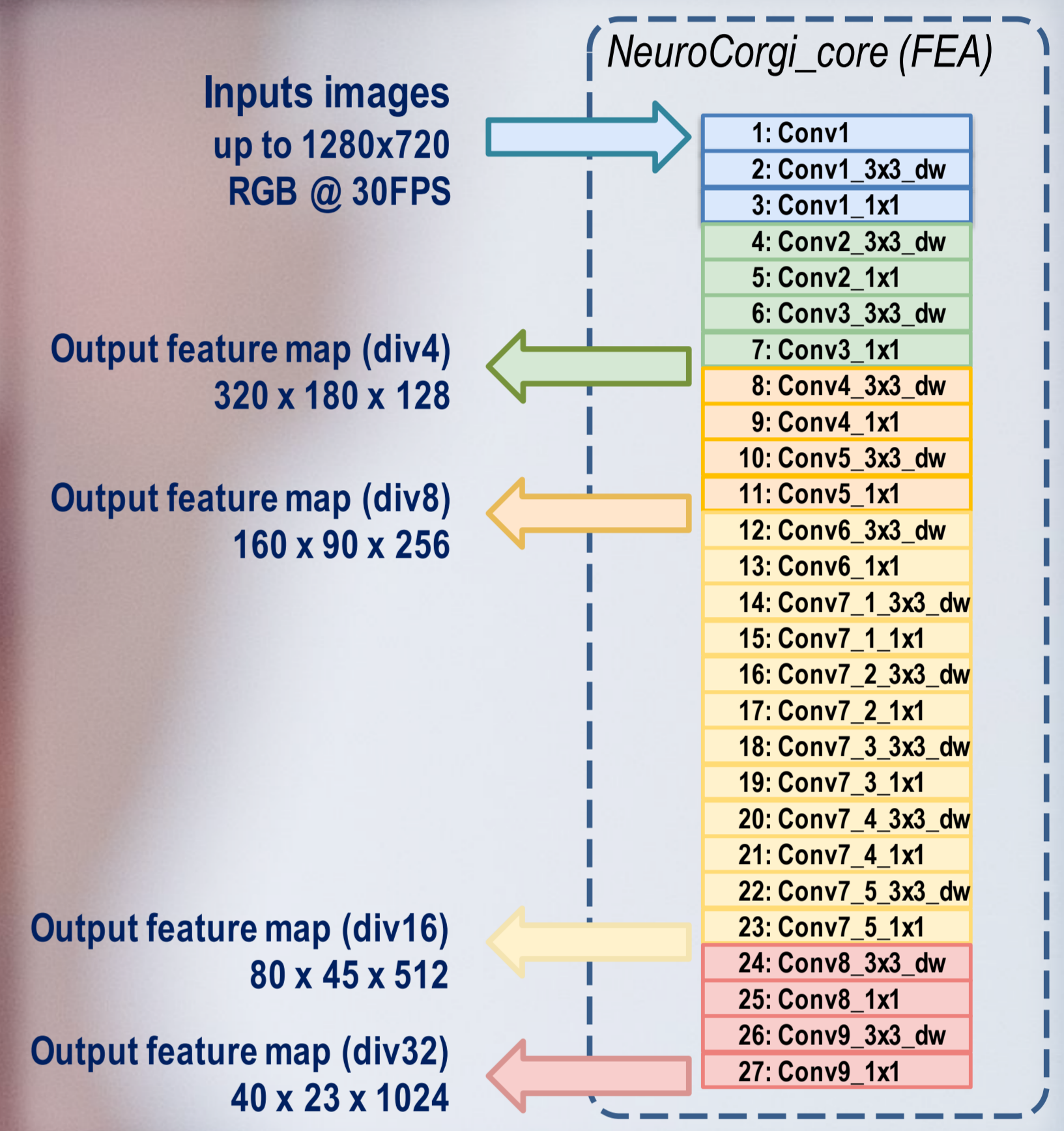


- NeuroCorgi is a Feature Extractor Accelerator (FEA) circuit
- Its Key design principle draws inspiration from the Human Brain, which shows that generic Feature Extractors can be used across a range of applications
- Leveraging a fixed MobileNet-V1 backbone and the Transfer Learning concept, its power dissipation beats the best State-of-the-Art accelerators by 9x to 1000x
- Built in 22FDX, it addresses Classification, Semantic Segmentation and Multi-Object Detection applications
- Designed using a set of tools for training the model (Aidge) and generating a complete new circuit (CorgiBuilder) for a given AI model

## Architecture



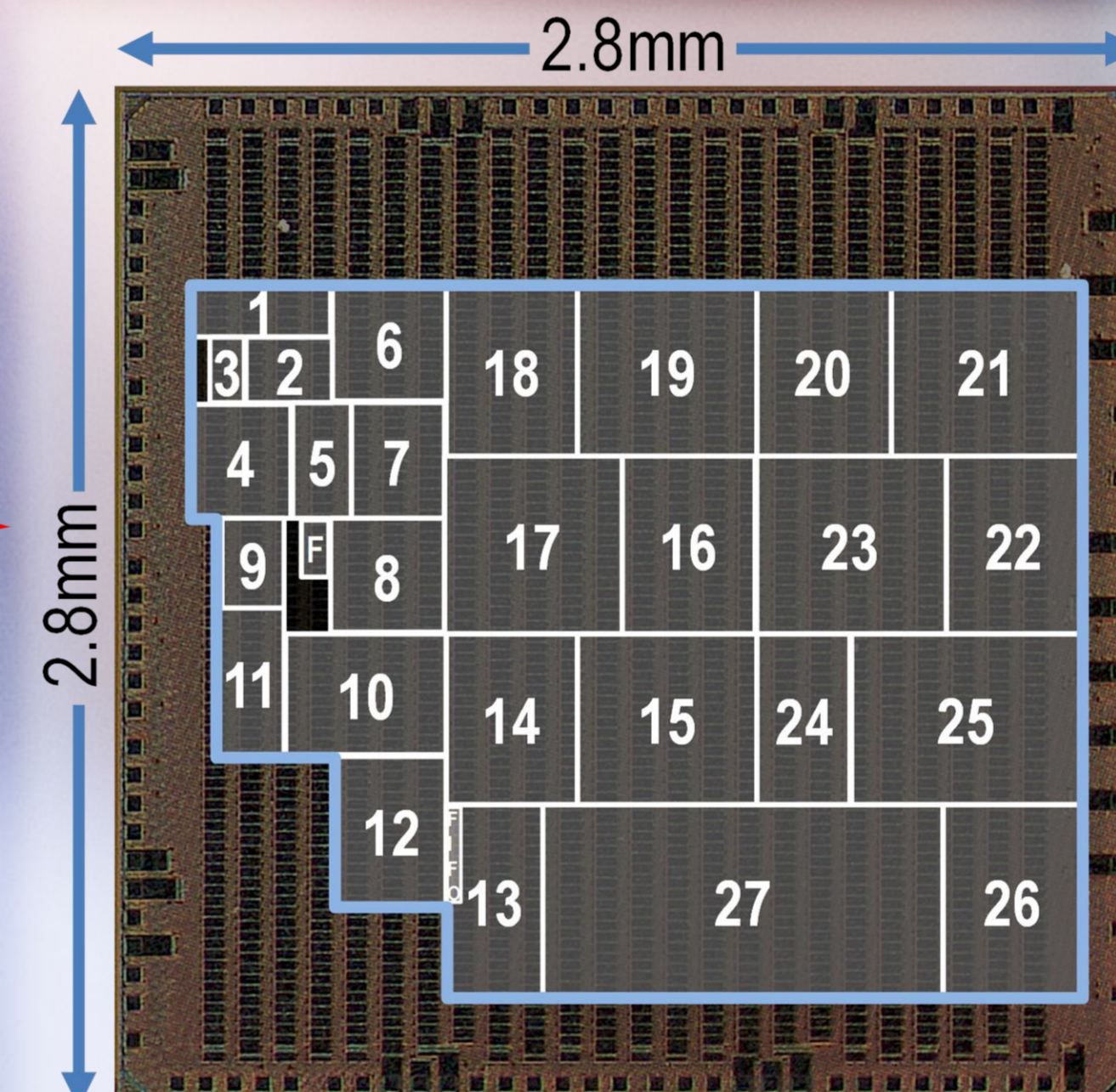
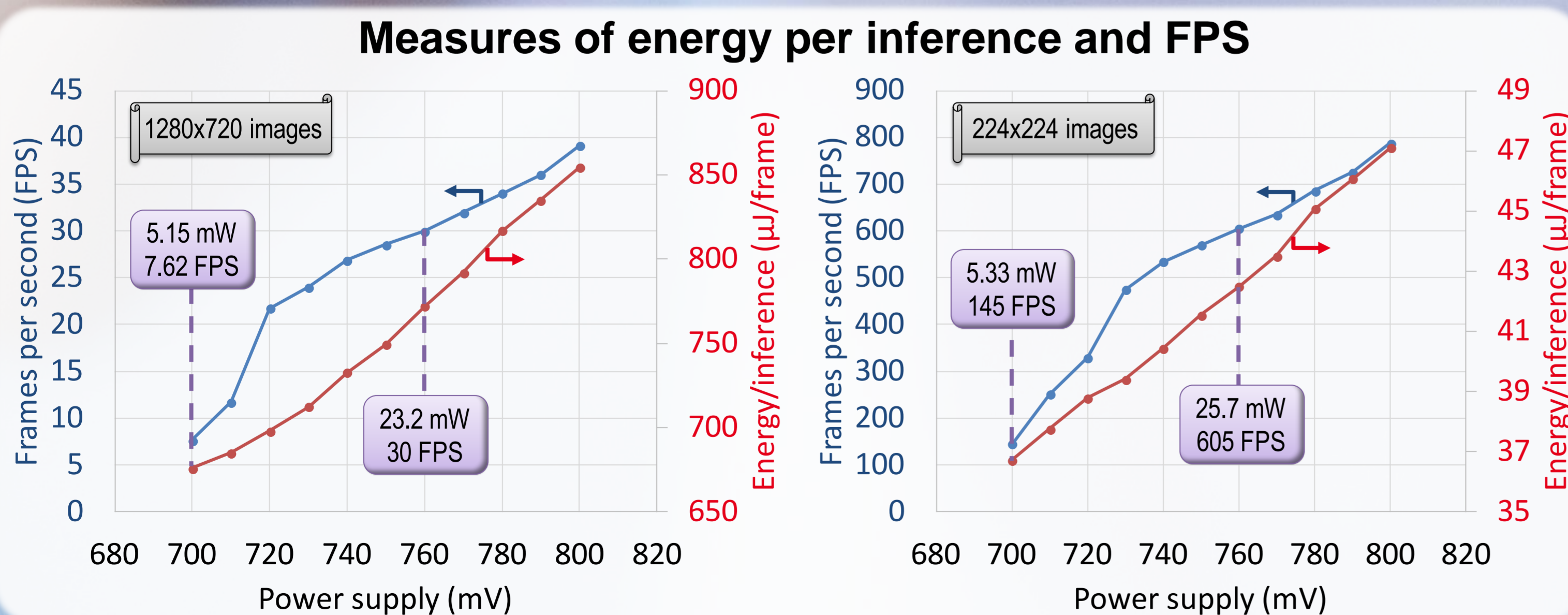
- Streaming architecture
- Support any image size up to HD (1280x720)
- MobileNet-V1 topology (27 layers)
- Trained on 4b weights and activations
- Accuracy of 70.42% on ImageNet
- 73.3 mIoU on Cityscapes semantic segmentation



## Technology

- NeuroCorgi is fabricated in GF 22FDX and comes in 3 variants
  - 'NeuroCorgi ImageNet' for Classification and Segmentation
  - 'NeuroCorgi NVM' for Classification with an OxRAM-based classifier
  - 'NeuroCorgi Coco' for Object Detection
- 'NeuroCorgi ImageNet' illustrated here contains
  - 3.2M parameters
  - 42kMAC operators
  - 186 SRAM memories

- 'NeuroCorgi ImageNet' measurements show outstanding results
- End-to-end energy efficiency: 30.9 TOPS/W (4b weights and activations)
- At least 9x better energy per inference than the State-of-the-Art
- Ultra-low power consumption:
  - 23.2mW for 1280x720 images at 30FPS
  - 1.5mW for 224x224 images at 30FPS



Chip summary	
Technology	GF 22FDX
Chip area	7.86mm <sup>2</sup>
FEA area	4.45mm <sup>2</sup>
# multipliers	42k
# SRAM memories	186
SRAM memory	1.1MB
Main clock	59MHz
AI model	MobileNet v1
Training dataset	ImageNet
Batch size	1

## Results

	ISSCC'20 [1]	JSSC'23 [2]	JSSC'23 [3] DIANA	JSSC'24 [4] Marsellus	JSSC'24 [5] DynaPlasia	ISSCC'22 [6] Hiddenite	This work [7] NeuroCorgi	
Technology	12nm	4nm	22nm	22nm FDX	28nm	40nm	22nm FDX	
Application	Server	Mobile	Edge	AI-IoT	Embedded	Embedded	Embedded	
Area (mm <sup>2</sup> )	709 (chip)	4.74 (core)	3.3 (core)	18.7 (chip) 2.42 (core)	20.25 (chip)	9 (chip) 4.36 (core)	7.86 (chip) 4.45 (FEA)	
Programmability	Programmable	Programmable	Programmable	Programmable	Programmable	Configurable	Fixed	
Power consumption (mW)	25W - 276W	381 - 5133	-	12.8 - 123	261	85.4 - 534.7	5 - 37	
ImageNet use case	Training dataset							
	AI model	ImageNet	ImageNet	ImageNet	ImageNet	ImageNet	ImageNet	
	Precision (bits)	ResNet50 v1	MobileNet TPU	ResNet18	ResNet18	ResNet18	ResNet50	MobileNet v1
	Top-1 accuracy (%)	8	8	Analog + digital	RBE 4x 4b	9w, 8a	ternary (w), 8a	4w, 4a
	Inferences/second (FPS)	74.93	-	64.1	68.5	70.4	70.09	70.42
		224x224	78563	3433	277	20.8	776 <sup>6</sup>	169.7 <sup>7a</sup>
	FEA latency (ms)	1280x720	-	-	-	-	-	-
		224x224	0.2 <sup>8b</sup>	0.29 <sup>9</sup>	3.6 <sup>9</sup>	48 <sup>9</sup>	1.29 <sup>6b</sup>	5.92 <sup>6b</sup>
	TOPS/W	1280x720	-	-	-	-	-	6.90
		224x224	4.14	11.59	5.52 <sup>2c</sup>	5.83	10.8 <sup>6</sup>	30.9 <sup>6</sup>
Best FEA Energy/inference (μJ/frame)	1280x720	-	-	-	-	-	30.9 <sup>6</sup>	
	224x224	2000 <sup>6</sup>	340 <sup>6</sup>	659 <sup>6b</sup>	557 <sup>6</sup>	336 <sup>6b</sup>	503 <sup>6b</sup>	
	1280x720	-	-	-	-	-	676	

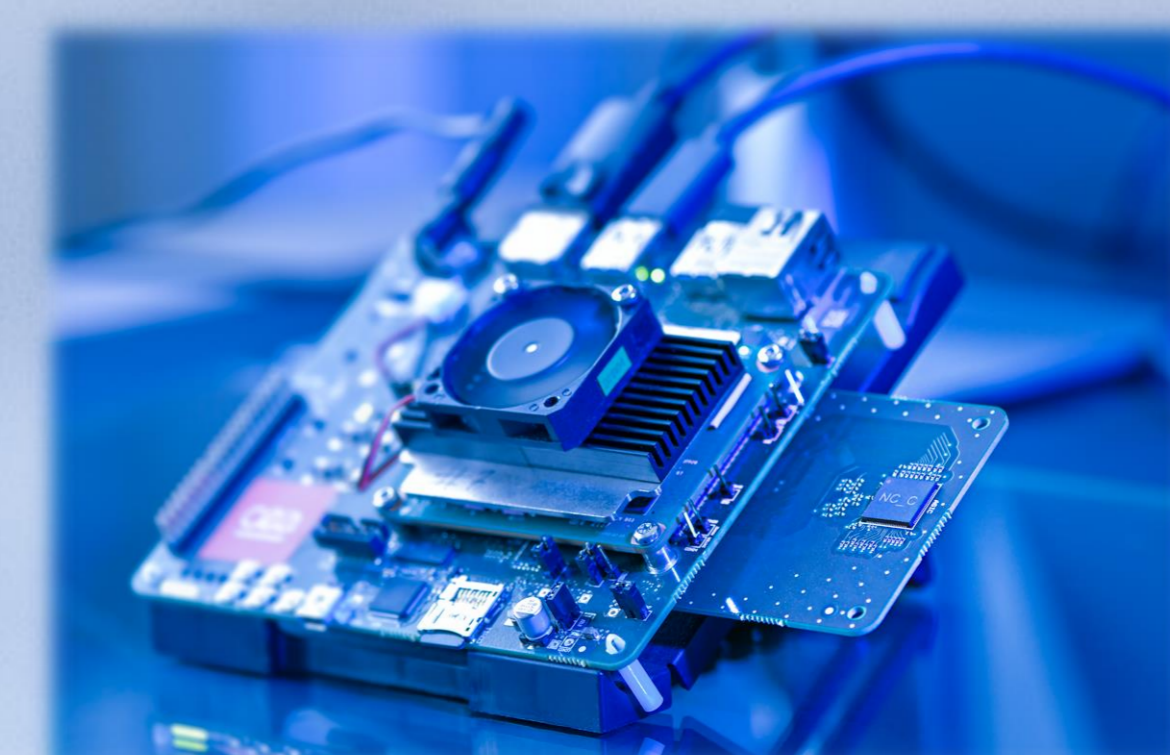
<sup>1</sup>MAC = 2 Ops. Zero skipping included as MACs <sup>4</sup>Without considering off-chip memory accesses. <sup>5</sup>Latency reported on Inception v3. <sup>6</sup>Only feature extraction (no FC layer) <sup>7</sup>Estimated feature extraction part. Assuming 1Conv OP = 1FC OP for latency and energy. FC is 0.18% (0.03%) of total MACs/frame on 224x224 images for MobileNet v1 (ResNet18) <sup>8</sup>Power & latency off-chip weight loading from external DRAM / external host CPU / on-chip network / on-chip memory access / refresh are not included.

## Acknowledgments



Reference:

- Y. Jiao et al., International Solid-State Circuits Conference, 2020.
- J. -S. Park et al., Journal of Solid-State Circuits, 2023.
- P. Houshmand et al., Journal of Solid-State Circuits, 2023.
- F. Conti et al, Journal of Solid-State Circuits, 2024.
- S. Kim et al., Journal of Solid-State Circuits, 2024.
- K. Hirose et al., International Solid-State Circuits Conference, 2022.
- I. Miro-Panades et al., Asia Pacific Conference on Circuit and Systems, 2024.



Scan Me to visit website