# How can efficient, low-latency and high-accuracy inference be performed in ADAS?
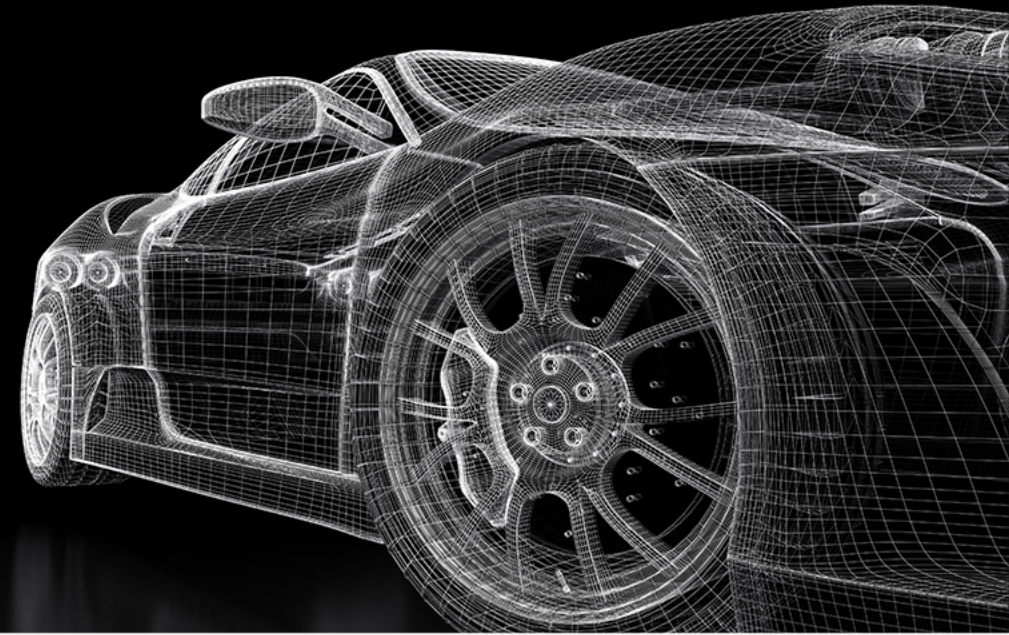
Kristofor Carlson, PhD

Manager of Applied Research

BrainChip Inc.

Another key efficiency feature of the VISION EQXX that takes its cue from nature is the way it thinks. It uses an innovative form of information processing called **neuromorphic computing.** The hardware runs spiking neural networks. Information is coded in discrete spikes and energy is only consumed when a spike occurs, which reduces energy consumption by orders of magnitude. **Working with California-based artificial intelligence experts BrainChip, Mercedes-Benz engineers developed systems based on BrainChip's Akida hardware and software.** The example in the VISION EQXX is the "Hey Mercedes" key-word detection. **Structured along neuromorphic principles, it is five to ten times more efficient than conventional voice control.**

Although neuromorphic computing is still in its infancy, systems like these will be available on the market in just a few years. When applied on scale throughout a vehicle, they have the potential to **radically reduce the energy needed to run the latest AI technologies.**

## In-Cabin AI

- Visual driver authentication
- Keyword spotting
- Voice authentication
- Contextual understanding

BrainChip is revolutionizing the future of in-device Artificial Intelligence (AI) and is the worlds first commercial producer of neuromorphic semiconductor chips and IP.

Our technology brings commonsense to the processing of sensor data, freeing machines to do more with less. Accurately. Elegantly. Meaningfully. We call this Essential AI.

Essential is optimizing compute. Maximizing performance. Minimizing power. In the real world. And in real time. We're proving that on chip AI, close to the sensor, has a sensational future, for our customers' products, as well as the planet.

**BrainChip. Essential AI.**

### BrainChip Profile:

- 15 years of AI architecture research.

- World leading team of neuromorphic experts.

- Centers of engineering excellence in Australia, USA, France and India.
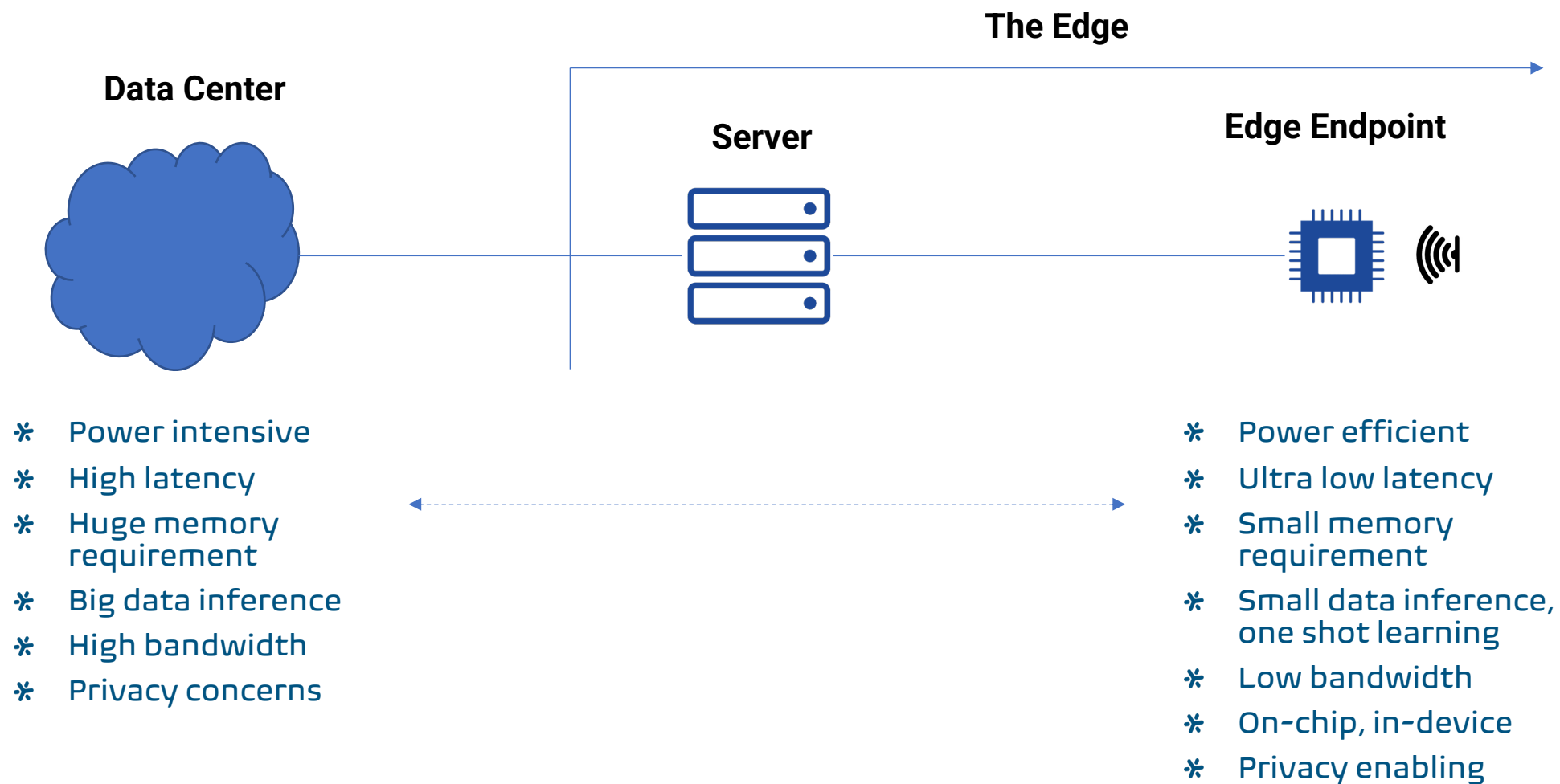
**Trusted By:**

MegaChips  Mercedes-Benz  NASA  RENESAS  Valeo

**Profiled In:**

VentureBeat  EE|Times  Design & Reuse

**Traded On:**

ASX AUSTRALIAN SECURITIES EXCHANGE  OTCQX

**The Edge**

**Data Center**

**Server**

**Edge Endpoint**

* Power intensive
* High latency
* Huge memory requirement
* Big data inference
* High bandwidth
* Privacy concerns

* Power efficient
* Ultra low latency
* Small memory requirement
* Small data inference, one shot learning
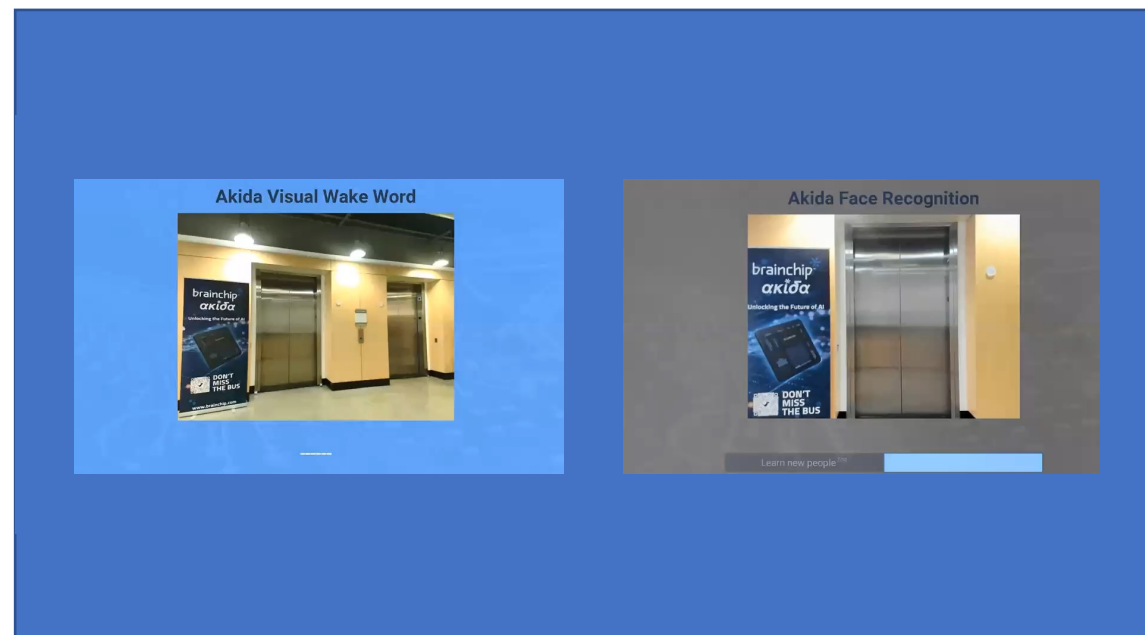* Low bandwidth
* On-chip, in-device
* Privacy enabling

## Vehicle and Person Detection



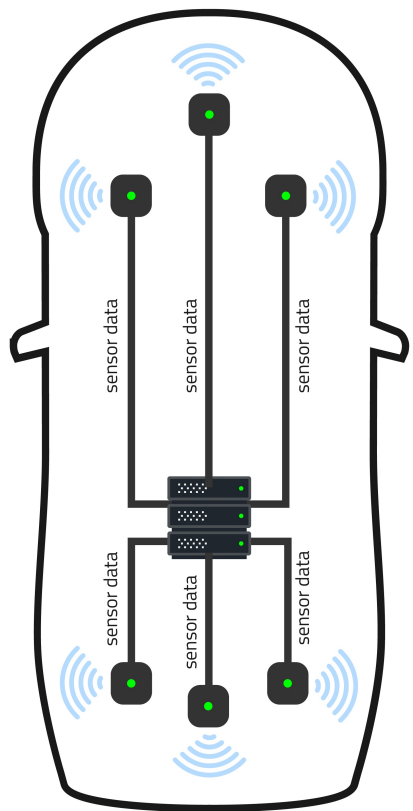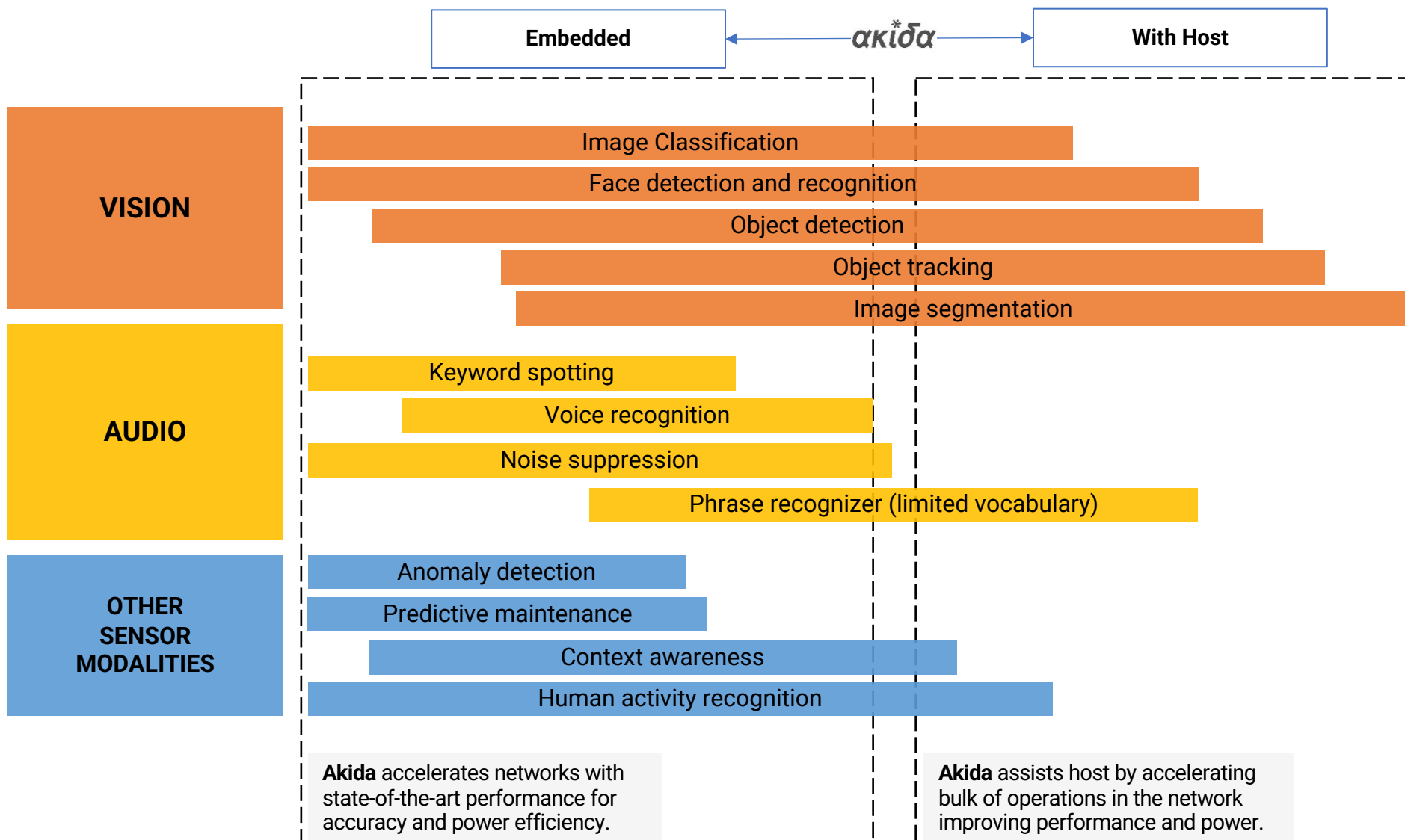## Visual Wake and Facial Recognition

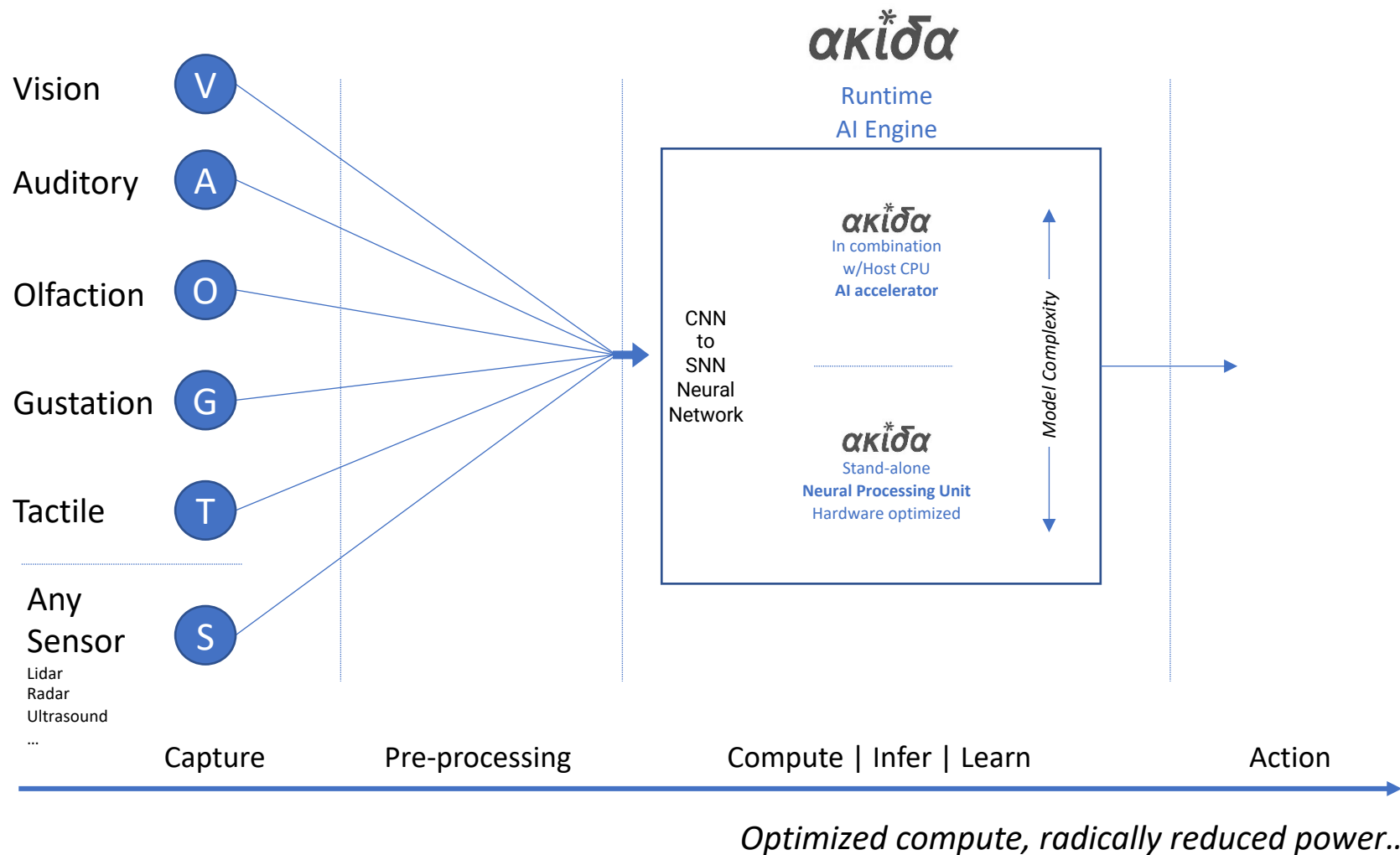**Conventional** Computing     VS     **On-Sensor** Computing

# AI Models to Support Key AI Use Cases

**Embedded** ⟷ ακἴδα ⟷ **With Host**

**VISION**
- Image Classification
- Face detection and recognition
- Object detection
- Object tracking
- Image segmentation

**AUDIO**
- Keyword spotting
- Voice recognition
- Noise suppression
- Phrase recognizer (limited vocabulary)

**OTHER SENSOR MODALITIES**
- Anomaly detection
- Predictive maintenance
- Context awareness
- Human activity recognition

**Akida** accelerates networks with state-of-the-art performance for accuracy and power efficiency.

**Akida** assists host by accelerating bulk of operations in the network improving performance and power.

## Performant and Efficient

Vision

Auditory

Olfaction

Gustation

Tactile

Any Sensor
Lidar
Radar
Ultrasound
…

**αкιδα**

Runtime
AI Engine

**αкιδα**
In combination
w/Host CPU
**AI accelerator**

CNN
to
SNN
Neural
Network

*Model Complexity*

**αкιδα**
Stand-alone
**Neural Processing Unit**
Hardware optimized

Capture    Pre-processing    Compute | Infer | Learn    Action

*Optimized compute, radically reduced power…*

- Akida supports output from any sensor, with its data and network, and applies in-chip neuromorphic AI to efficiently process and infer with radically reduced power consumption.

- Akida can also be deployed as an AI Accelerator to work in combination with a host CPU, still providing power consumption efficiencies.

**Power Efficient**
*Microwatts*

**Ultra Low
Latency**

**Low Memory
Footprint**

**One Shot
Learning**

**Flexible IP &
Quick to Deploy**

**Independent**
*Of the Cloud*

**Privacy Enabling**
*On-Device
Processing*

**Environment**
*Less Power,
Less Carbon*

**Distributed Computation**

*Each NPU has dedicated compute and memory, reducing data movement.*



**Event-Based Processing**

*NPUs perform computation only on events (non-zero values).*



**Event-Based Communication**

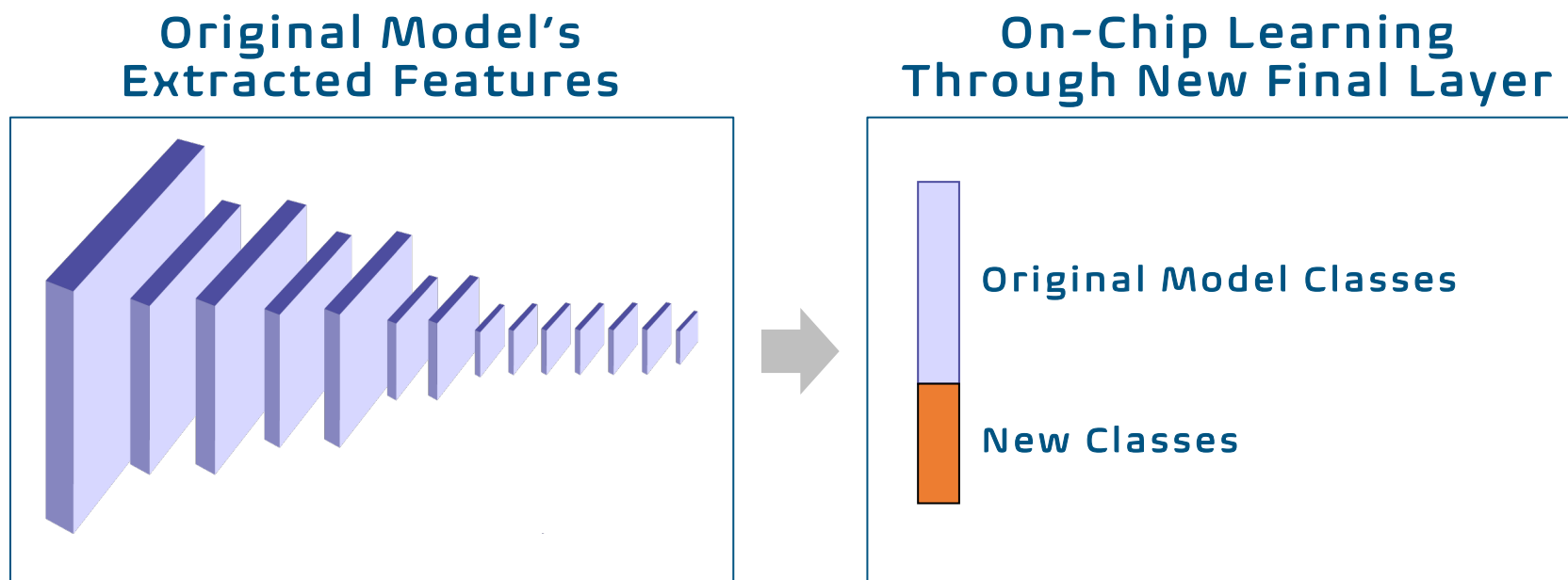*Send events over mesh network without host CPU intermediation.*



**Event-Based Learning**

*On-chip learning algorithm.*

* Brainchip IP and Akida chips are capable of on-chip learning by leveraging the trained model as a feature extractor then adding new classes to the final layer.
* Demonstrated Edge Learning for:
  - Object detection using MobileNet trained on the ImageNet dataset.
  - Keyword spotting using DS-CNN trained on the Google Speech Commands dataset.
  - Hand gesture classification using small CNN trained on a custom DVS events dataset.



**Original Model's Extracted Features**

**On-Chip Learning Through New Final Layer**

Original Model Classes

New Classes

✱ Multi-pass is a component of the Brainchip architecture that reduces the number of neural processing units required for a given compute task by segmenting and processing sequentially.
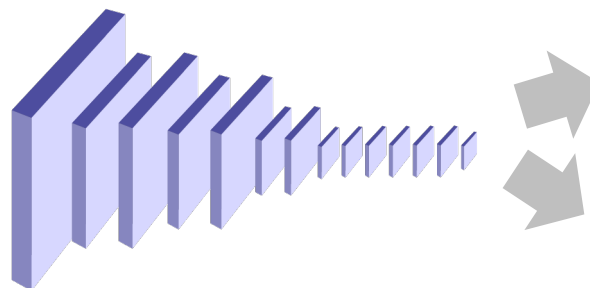
## Benefits of Multi-Pass

## How it Works

**Scalable**

**Smaller Memory Requirement** *(2X)*

**Power Efficient**

**Example CNN**

**Multi-pass Sequential Compute**

Node    Repeat    .....    .....

**Parallel Compute in One-Step**

Node    Node    1 Step
Node    Node

| Task | Visual Wake Words | Image Classification | Keyword Spotting | Anomaly Detection |
|---|---|---|---|---|
| Data | VWW | CIFAR-10 | Google Speech | ToyADMOS |
| Model | MobileNetV1 (0.25x) | Resnet-V1 | DS-CNN | FC AutoEncoder |
| Harvard Arm Cortex-M4 w/FPU* | 603ms / 24,320μJ/inf | 704ms / 29,207μJ/inf | 181ms / 7,373μJ/inf | 10ms / 416μJ/inf |
| PCL RV32IMAC w/FPU* | 846ms | 1,239ms | 325ms | 14ms |
| Syntiant Cortex-M0* | - | - | 6ms | - |
| Google Coral** | - | - | 0.5ms / 351μJ/inf | - |
| NVIDIA Nano** | - | - | 2.2ms / 659μJ/inf | - |

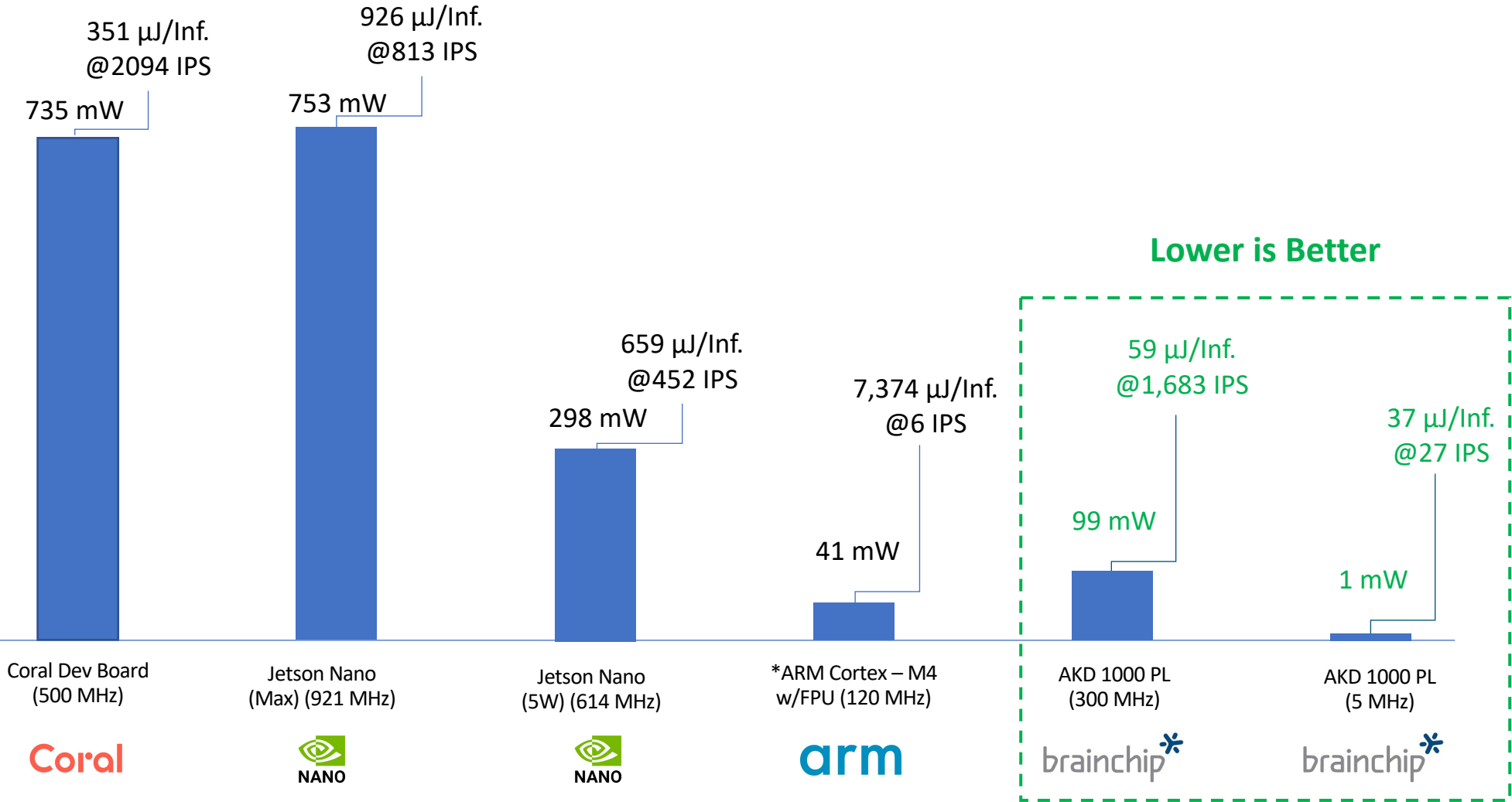| | | | | |
|---|---|---|---|---|
| Brainchip Akida1000** | 13ms / 259μJ/inf | | 3ms / 66μJ/inf | - / 4μJ/inf |
| **Brainchip Improvement Over Best Performing Chip*** | **65x Faster 94x Less Power** | | **Similar Speed 5-10x Less Power** | **104x Less Power** |

*MLPerf™ v0.5 Inference Closed ResNet-v1.5 offline. Retrieved from https://mlcommons.org/en/inference-tiny-05/ 20 April 2022, entries 0.5-464, 0.5-465, and 0.5-468. The MLPerf name and logo are trademarks of MLCommons Association in the United States and other countries. All rights reserved. Unauthorized use strictly prohibited. See www.mlcommons.org for more information.

**Result not verified by MLCommons.

***Brainchip Improvement denotes the Brainchip Akida1000 performance compared to the next best performing chip in each category given available data.

351 μJ/Inf. @2094 IPS — 735 mW — Coral Dev Board (500 MHz)

926 μJ/Inf. @813 IPS — 753 mW — Jetson Nano (Max) (921 MHz)

659 μJ/Inf. @452 IPS — 298 mW — Jetson Nano (5W) (614 MHz)

7,374 μJ/Inf. @6 IPS — 41 mW — *ARM Cortex – M4 w/FPU (120 MHz)

Lower is Better

59 μJ/Inf. @1,683 IPS — 99 mW — AKD 1000 PL (300 MHz)

37 μJ/Inf. @27 IPS — 1 mW — AKD 1000 PL (5 MHz)

*MLPerf™ v0.5 Inference Closed ResNet-v1.5 offline. Retrieved from https://mlcommons.org/en/inference-tiny-05/ 20 April 2022, entries 0.5-464. The MLPerf name and logo are trademarks of MLCommons Association in the United States and other countries. All rights reserved. Unauthorized use strictly prohibited. See www.mlcommons.org for more information.
** All other results were not verified by MLCommons.

## Akida AKD1000 is a reference chip implemented with TSMC at 28nm which proved viability of IP.
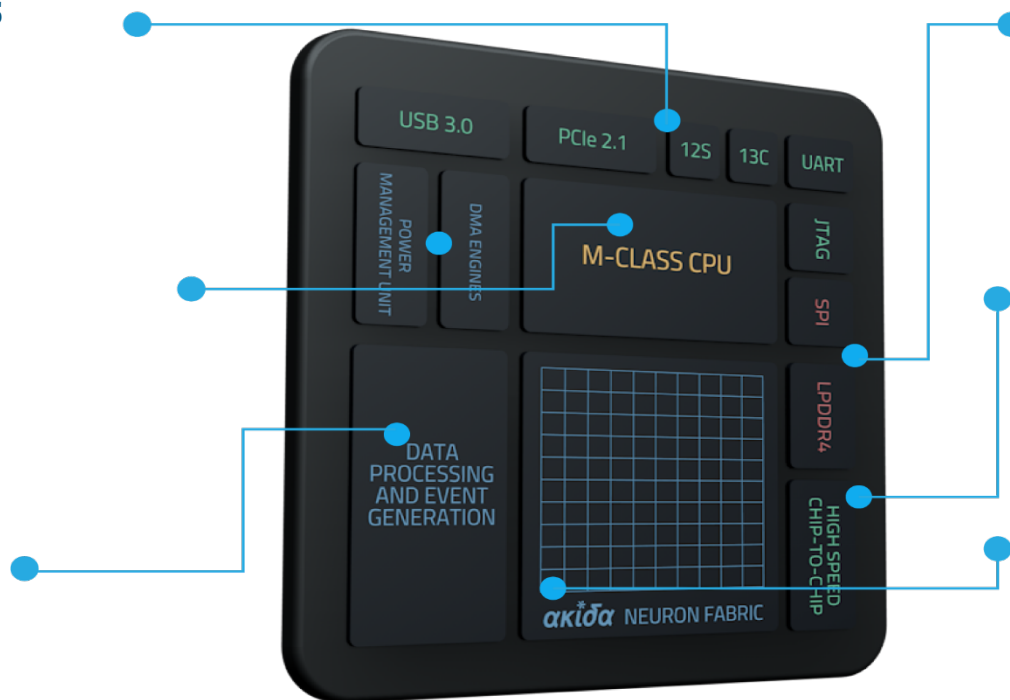
### Data Input Interfaces

- PCI Express 2.1 x2 Lane Endpoint
- USB 3.0 Endpoint
- I3S, I2C, UART, JTAG

### On-Chip Processor

- M-Class CPU with FPU & DSP
- System Management
- Akida Configuration

### Data Processing

- Pixel-Event Converter
- SW Data-Event Encoder
- Any multivariable digital data
- Sound, pressure, temp., others



### External Memory Interfaces

- SPI FLASH for boot/storage
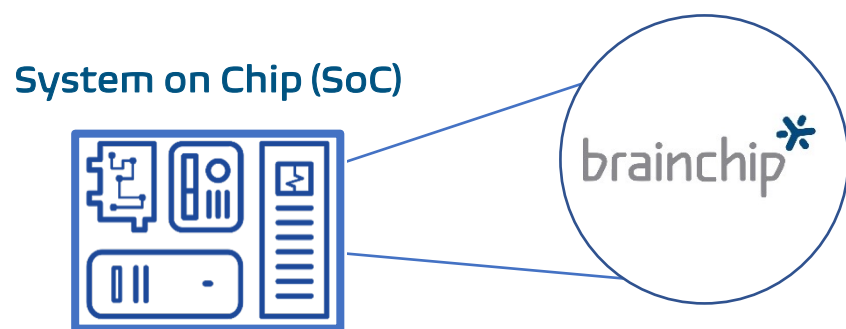- LPDDR4 Program/Weights

### Multi-Chip Expansion

- PCIe 2.1 2 lane root complex
- Connects up to 64 devices

### Flexible Akida Neuron Fabric

- Implements 80 NPUs
- All Digital logic with SRAM (8MB)
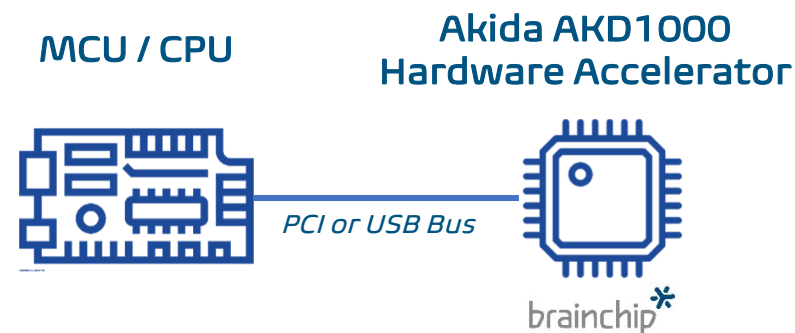- Also Available as Licensed IP Core
- First Implementation: TSMC 28nm

## IP Integrated into SoC

### System on Chip (SoC)



brainchip*

- Brainchip IP can be designed into a System on Chip (SoC) to seamlessly accelerate all AI workflows.
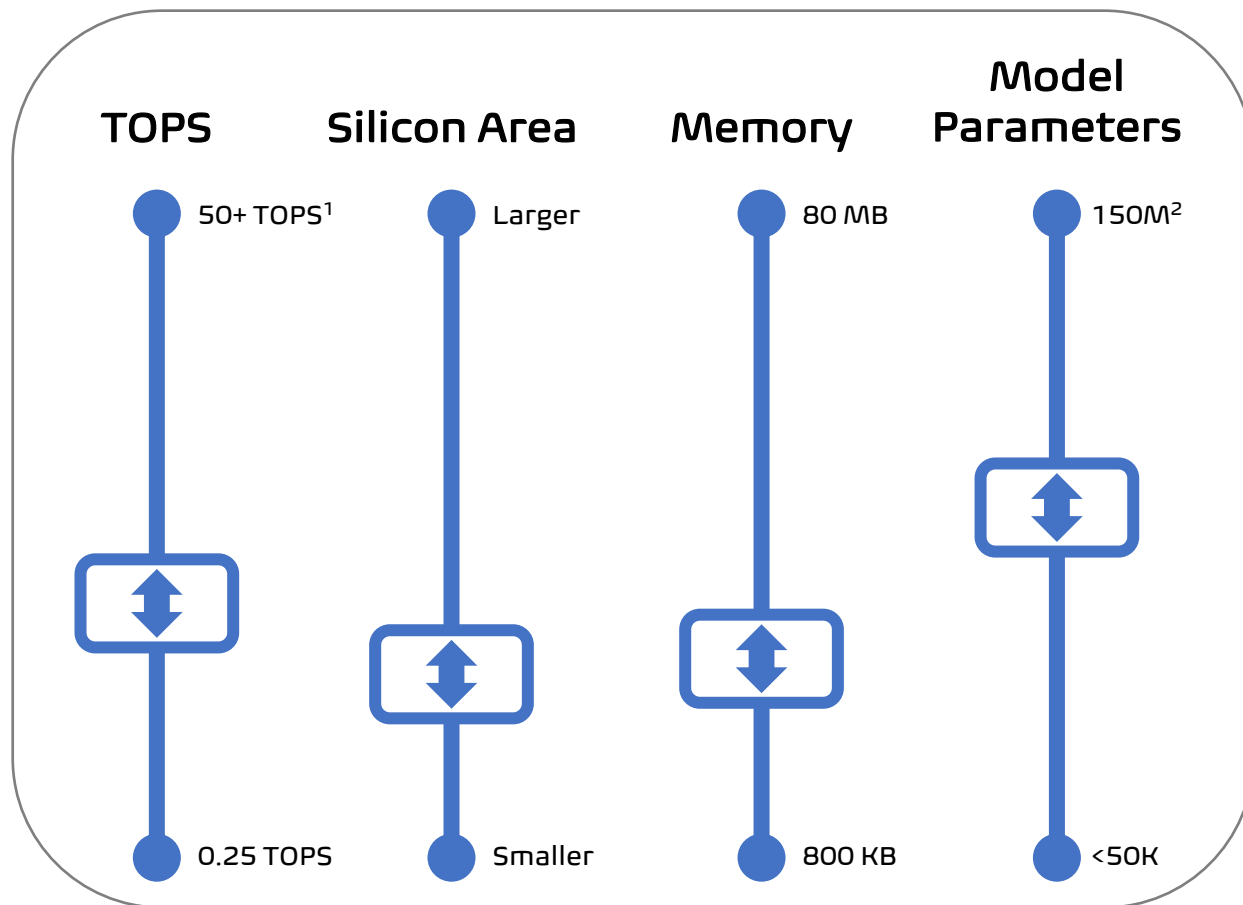- Applications can be ported to SoC and be adapted to market requirements; cost, size, and power.

## Reference Chip as a Hardware Accelerator

### MCU / CPU     Akida AKD1000 Hardware Accelerator



*PCI or USB Bus*

brainchip*

- ✳ An Akida AKD1000 with accompanying software development ecosystem can be placed next to an MCU or CPU where it can accelerate AI compute tasks.
- ✳ Brainchip's software development ecosystem and runtime is seamlessly compatible with any CPU; hardware and OS agnostic.

**BrainChip IP is flexible and scalable and can be implemented to support multiple edge AI use cases.**
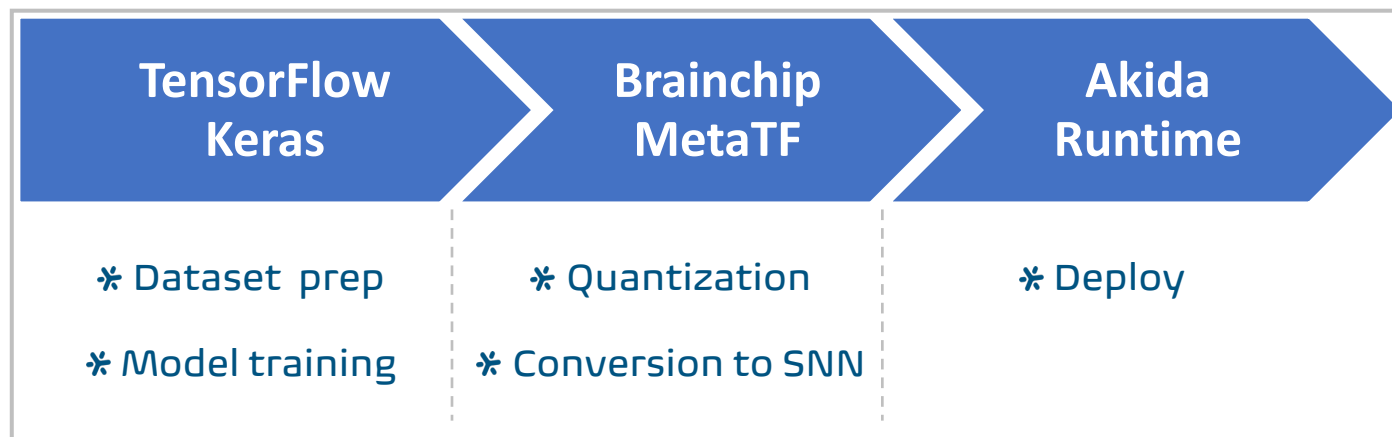
| TOPS | Silicon Area | Memory | Model Parameters |
|------|-------------|--------|-----------------|
| 50+ TOPS[1] | Larger | 80 MB | 150M[2] |
| 0.25 TOPS | Smaller | 800 KB | <50K |

- Brainchip works with clients to achieve the most cost-effective solution by optimizing the node configuration to the desired level of performance and efficiency.

- Scale down to 2 nodes for ultra low power or scale up to 256 nodes for complex use cases.

- Multi-pass processing provides flexibility to process complex use cases with fewer nodes increasing power efficiency.

- Quantization in MetaTF converts model weights and activations to lower bit format reducing memory requirement.

Notes:
1. 50+ TOPS is based on 100 Nodes at 1Ghz.
2. With optional SRAM and host CPU

**Our open-source software to easily convert models for native Akida runtime.**

- Create models in TensorFlow Keras.

- Easily convert TensorFlow Keras models with MetaTF API.
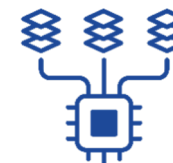
- No need to learn a new ML framework.

| TensorFlow Keras | Brainchip MetaTF | Akida Runtime |
|---|---|---|
| ✳ Dataset prep | ✳ Quantization | ✳ Deploy |
| ✳ Model training | ✳ Conversion to SNN | |

We make AI enablement easy…

**Define Use Cases**  →  **Create Data Set**  →  **Import with MetaTF**  →  **Optimize Network & Hardware Configuration**
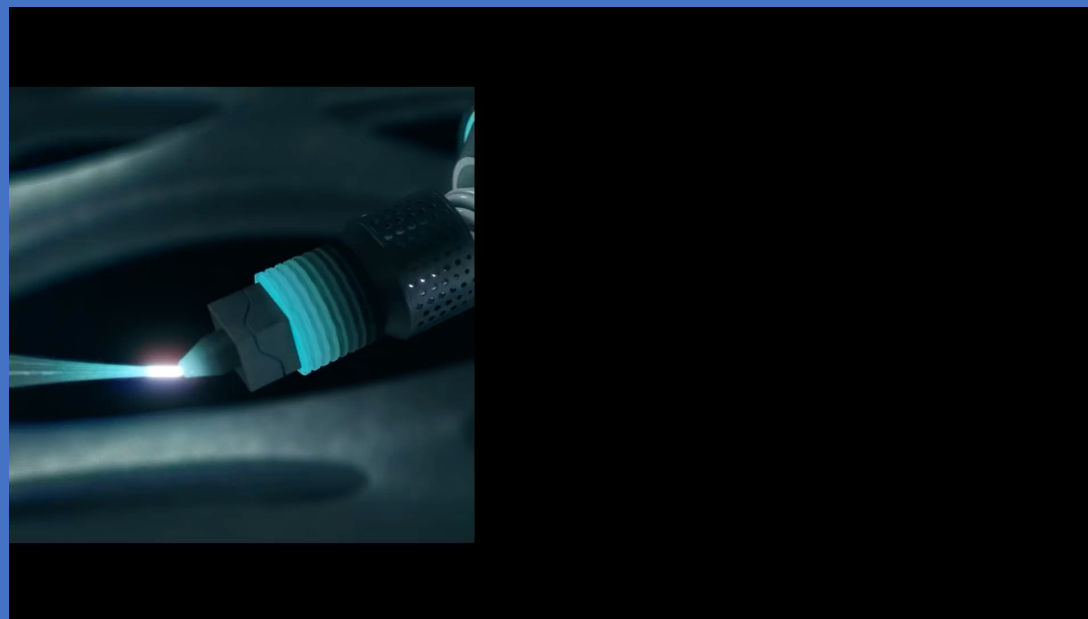
## Early Adopters

## Licensees

Most customers cannot be identified due to Non-Disclosure Agreements.

"Hey Mercedes"

# THANK YOU.
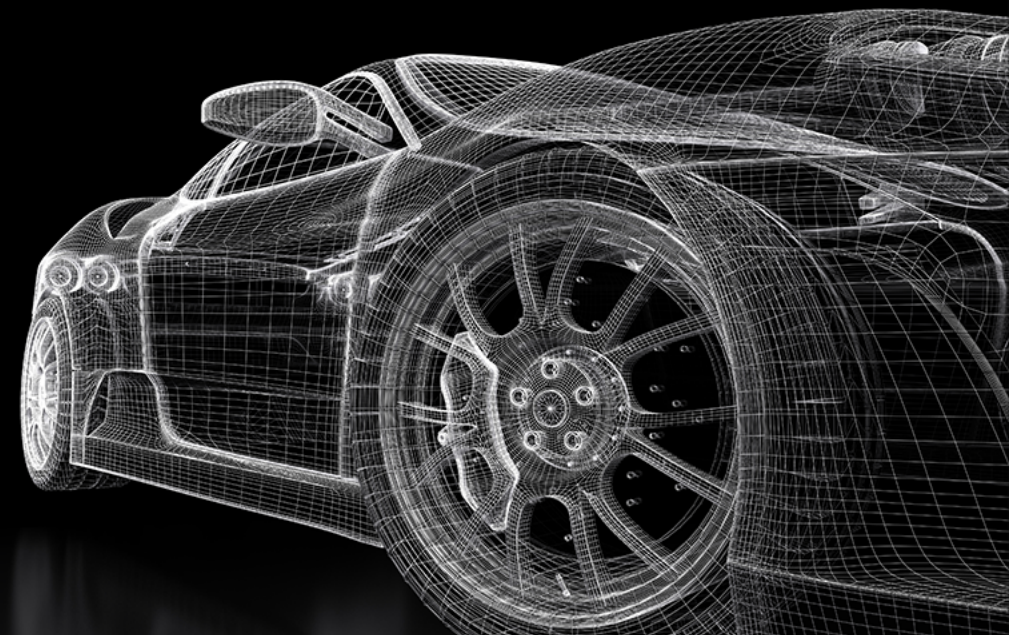
# QUESTIONS?

Kristofor Carlson, PhD

Manager of Applied Research

BrainChip Inc