



# Energy-efficient Semiconductors for AI in the age of Large Language Models

Dr. Sakyasingha Dasgupta

CEO & Founder  
EdgeCortix

EDGECORTIX™

# EDGECORTIX<sup>®</sup>

## Our Mission

*“To deliver cloud-level performance at the edge, with orders of magnitude better energy efficiency and processing speed, drastically reducing customer operating costs.”*

We are pioneering the future of the connected intelligent edge with a Software-driven **Edge Artificial Intelligence (AI) Platform**

Backed By



### Strategic Partnership with Renesas in AI Acceleration



Across Microcontrollers & Microprocessors

**2023 & beyond ..**  
Collaboration to streamline AI/ML Development.  
[AI/ML Press Release](#)

**Oct 2023**

Renesas invests in EdgeCortix as part of \$20M in additional funding.  
[Funding PR Press Release](#)

**Sept 2022**

Renesas announces RZ/V will use MERA in DRP-AI.  
[RZ/V Press Release](#)

**July 2022**

Collaboration on DRP-AI TVM compiler announced.  
[DRP-AI Press Release](#)

Since 2021

EdgeCortix is headquartered in Japan, with offices in USA, Singapore and India

Founded July 2019



Among Global 20\* (Pitchbook)

EdgeCortix Engineering Center  
Mushashi-Kosugi, Kawasaki-Shi

# ARTIFICIAL INTELLIGENCE

how machine learning will shape the next decade

## WIRED

MATT BURGESS

# AI demand is changing how we design chips

## GENERATIVE AI



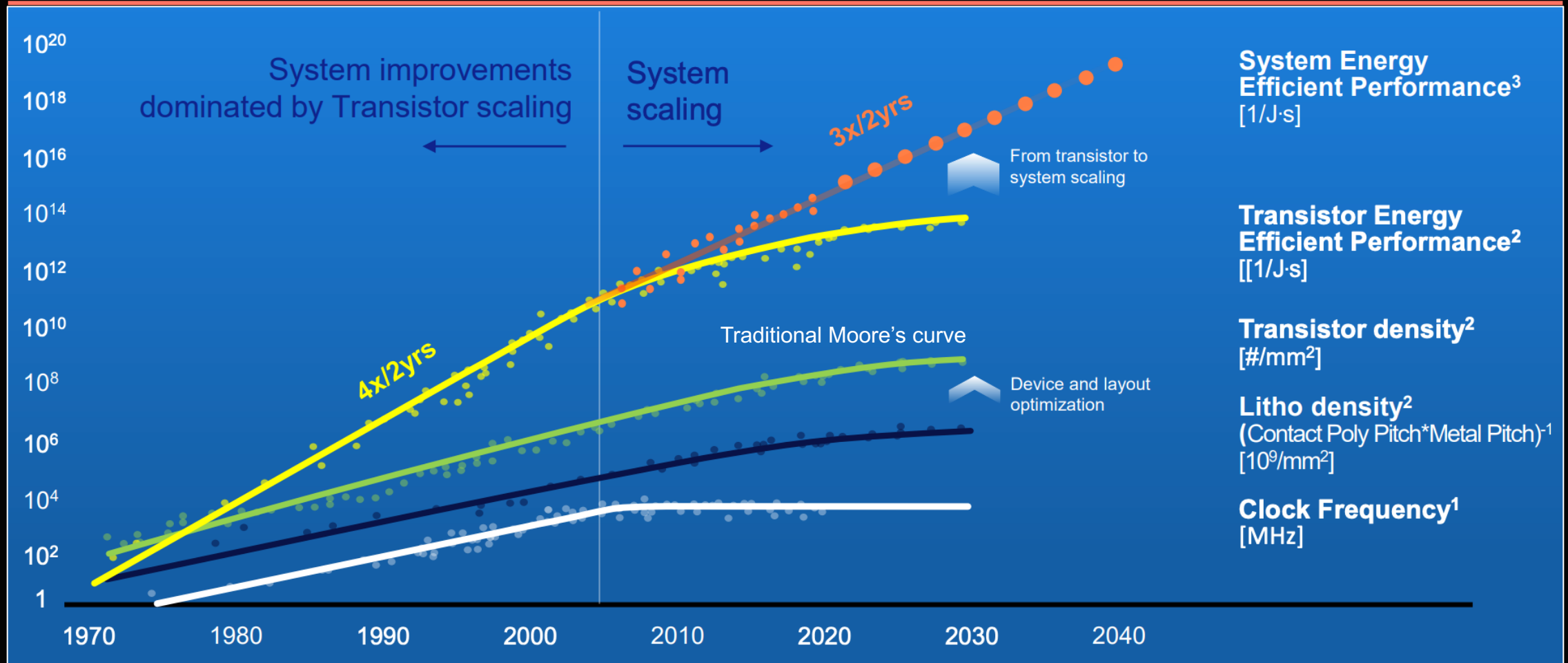
A NEW-AGE CREATOR AND DISRUPTOR



Image source: <https://www.purplequarter.com/generative-ai-new-age-creator-disruptor/all-about-tech/>

# Revisiting Moore's Law Evolution - An optimistic view

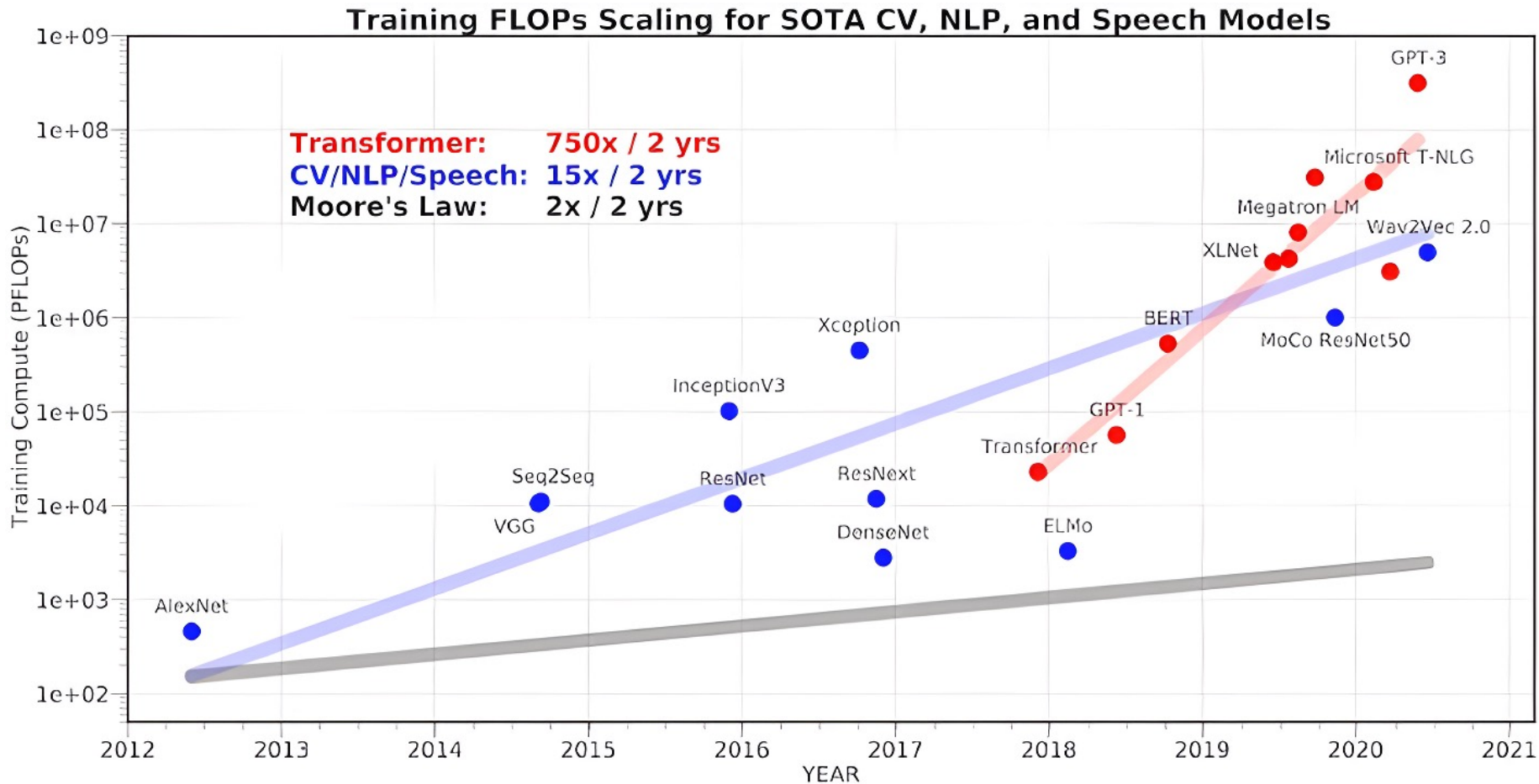
Moving into a System Scaling Era of Beyond Moore. Is this enough ?



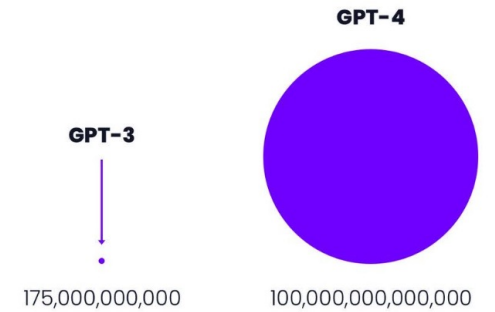
Sources: <sup>1</sup>Karl Rupp <sup>2</sup>ASML Presentation, <sup>3</sup>Mark Liu, TSMC, normalized to transistor EEP in 2005.

# Artificial Intelligence Computing Demand vs Moore's Law

Generative AI Driving a Hyper-exponential Demand in Computing: > 100x vs Moore's

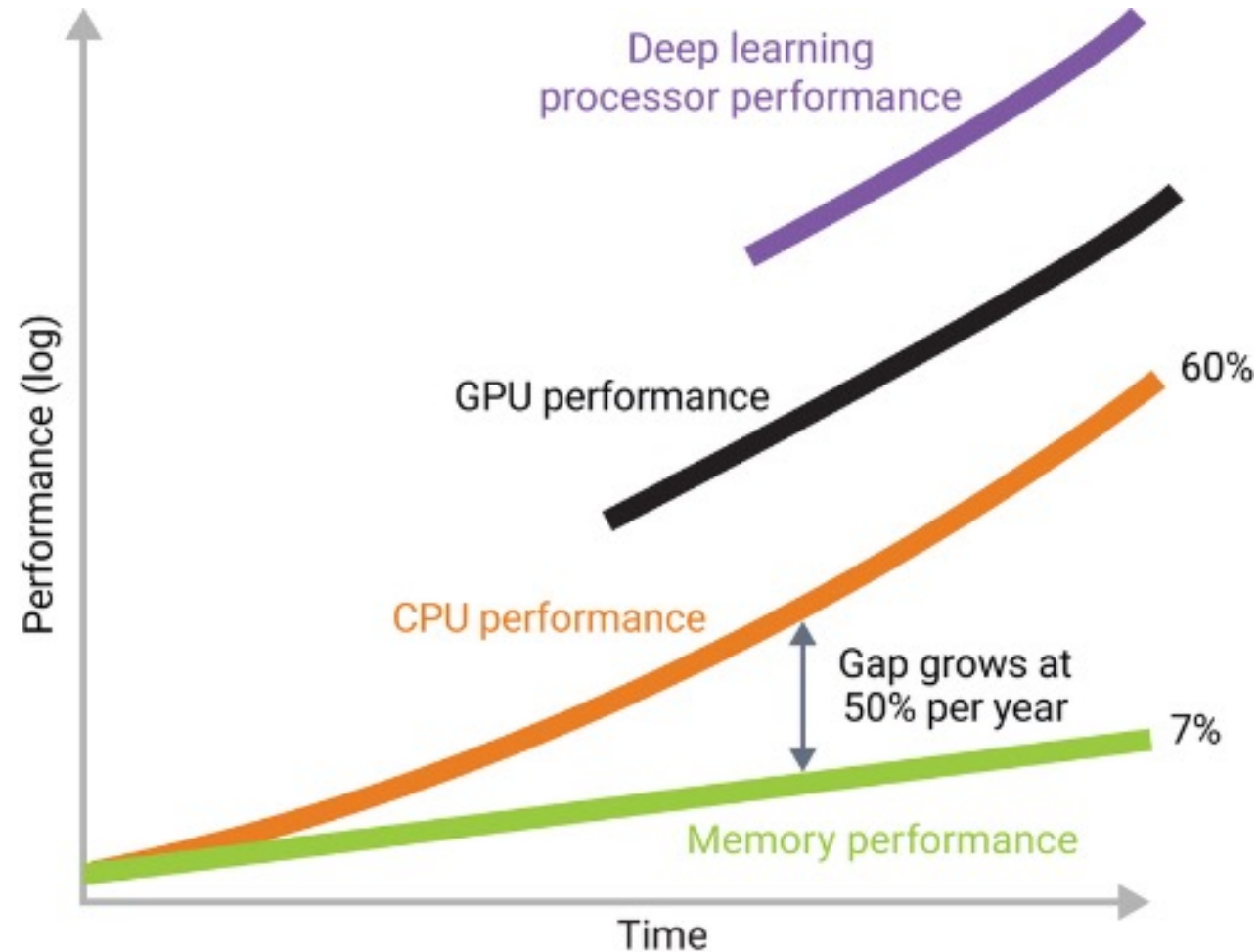


Recent Models like GPT-4 uses 1 Trillion Parameters or more ..



# Overcoming the Memory Wall of AI Processing

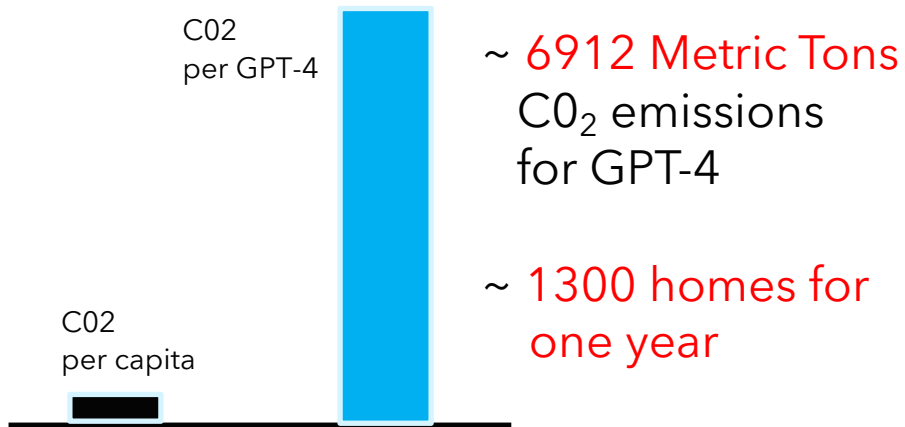
Memory is now a fundamental performance and energy bottleneck for AI Processing



# Not Just Performance, Need for Higher Energy-Efficiency

Need orders of magnitude improvements in AI processor performance/watt

## GPT-4 Carbon Footprint

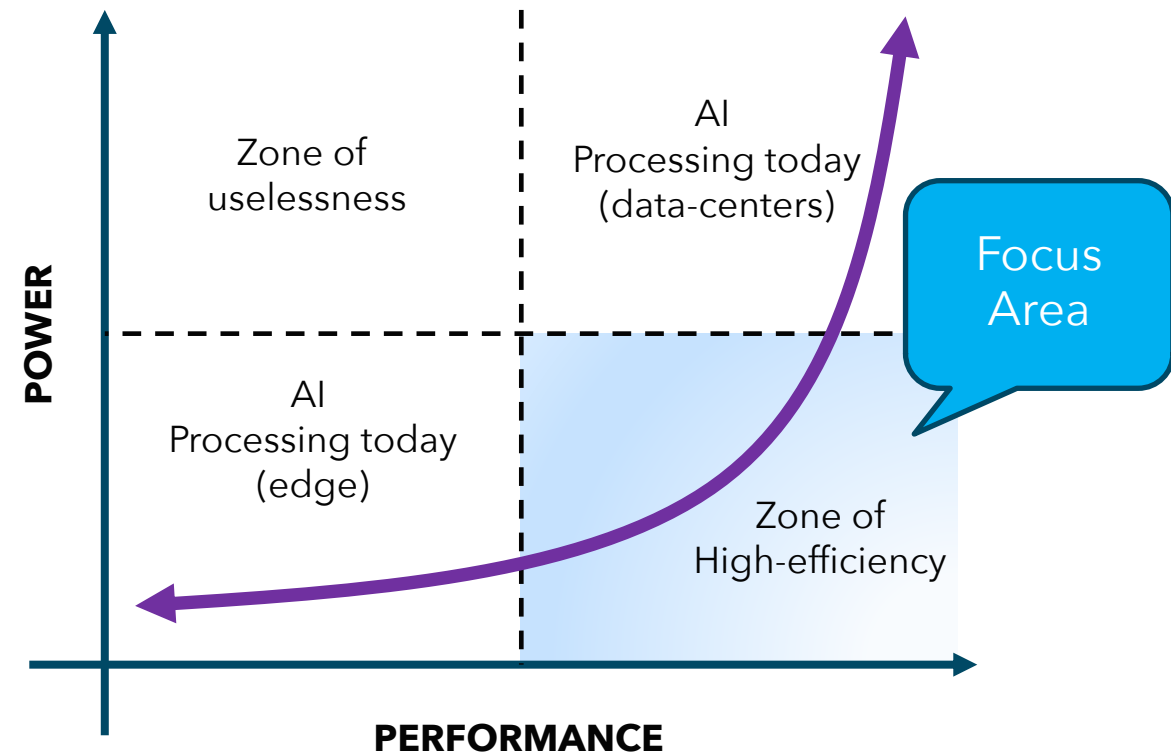


### Environmental Impact of Select Machine Learning Models, 2022

Source: Luccioni et al., 2022 | Table: 2023 AI Index Report

Model	Number of Parameters	Datacenter PUE	Grid Carbon Intensity	Power Consumption
Gopher	280B	1.08	330 gCO <sub>2</sub> eq/kWh	1,066 MWh
BLOOM	176B	1.20	57 gCO <sub>2</sub> eq/kWh	433 MWh
GPT-3	175B	1.10	429 gCO <sub>2</sub> eq/kWh	1,287 MWh
OPT	175B	1.09	231 gCO <sub>2</sub> eq/kWh	324 MWh

## Semiconductor Design Goal



# Enable High Performance, Low Power Processing at the Edge

Reduce Expensive Data Movement by Moving more AI Processing to the Site of Data Creation

**74 ZB**

Data generated at the edge<sup>1</sup>

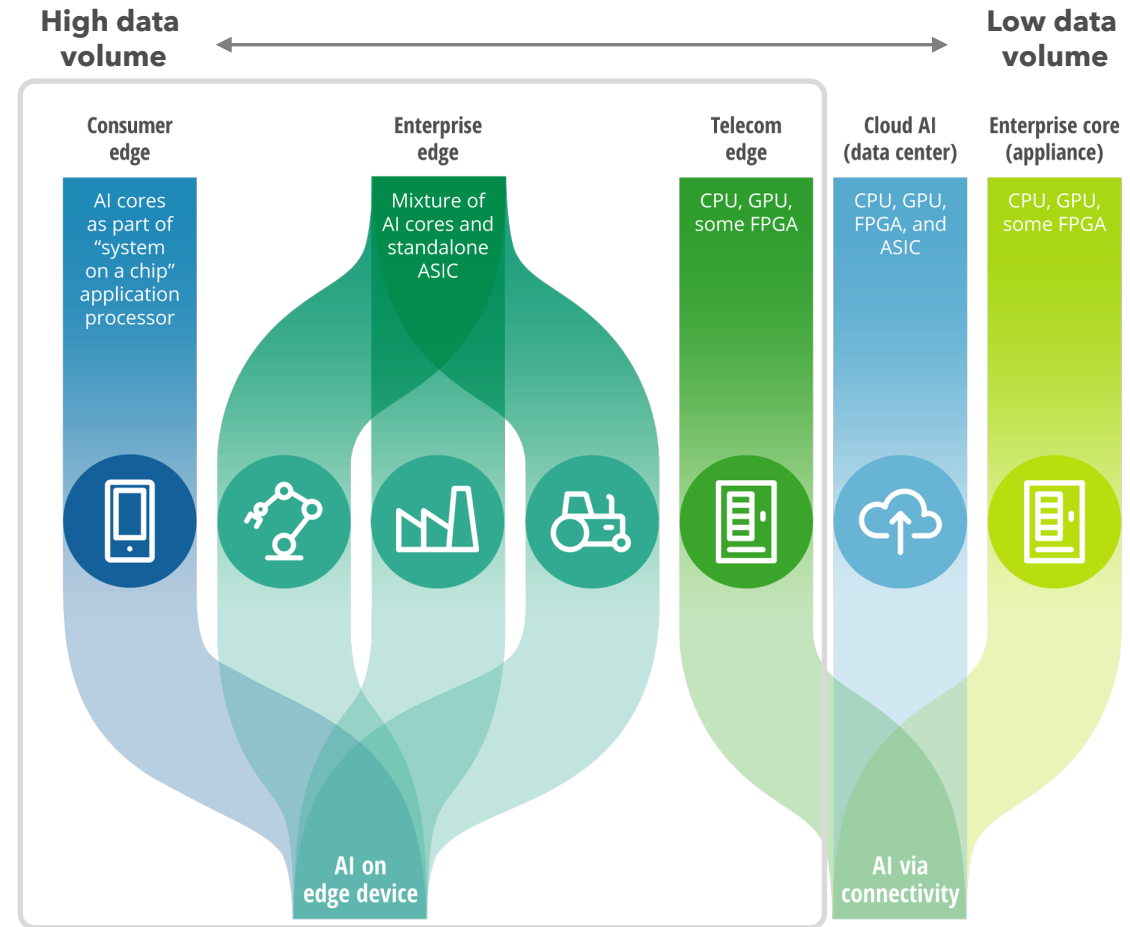
**75%**

Of enterprise-generated data created and processed at the edge by 2025<sup>2</sup>

**+\$80B**

Edge AI market<sup>3</sup>

Sources: <sup>1</sup>IDC annual estimate for 2025 <sup>2</sup>Gartner <sup>3</sup>IDC

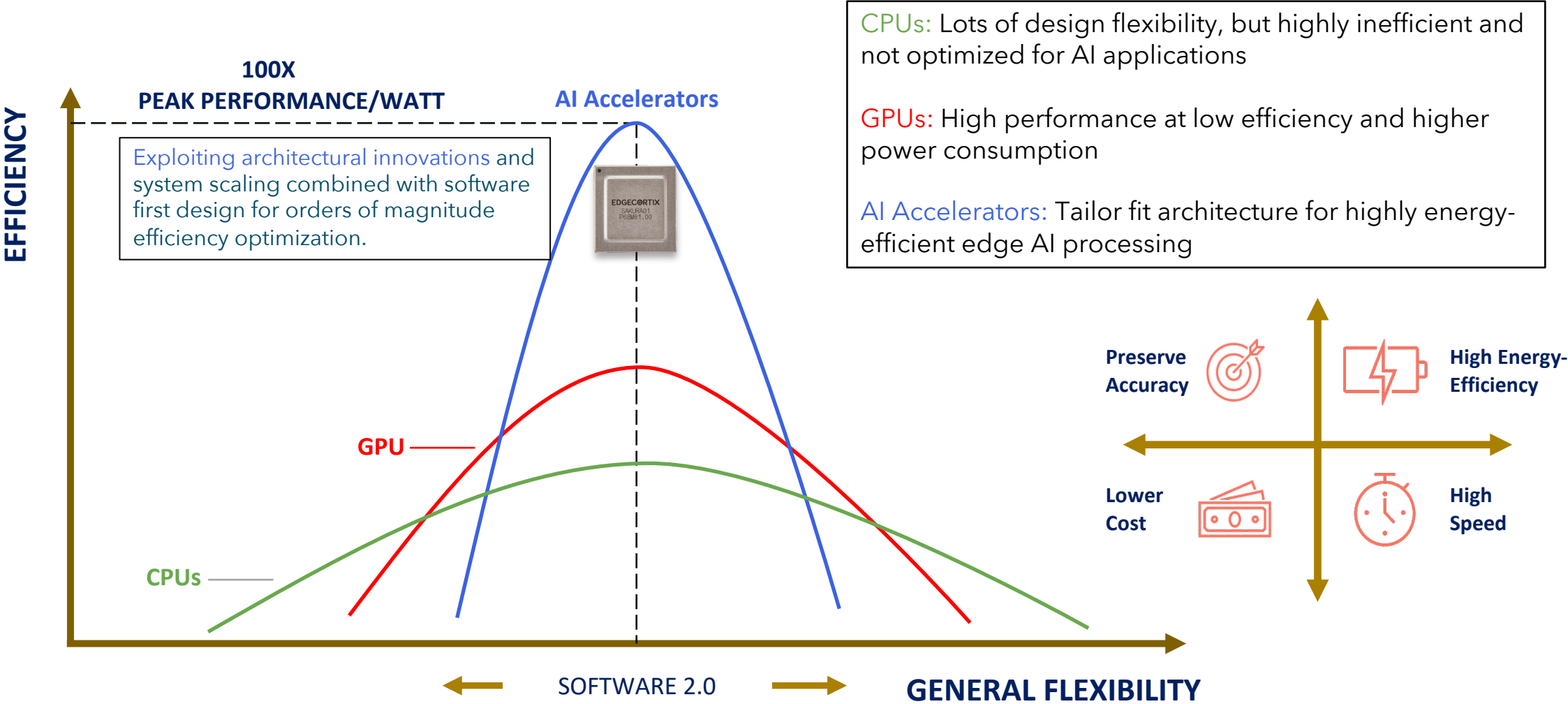


Source: Deloitte analysis.

Deloitte Insights | [deloitte.com/insights](https://deloitte.com/insights)

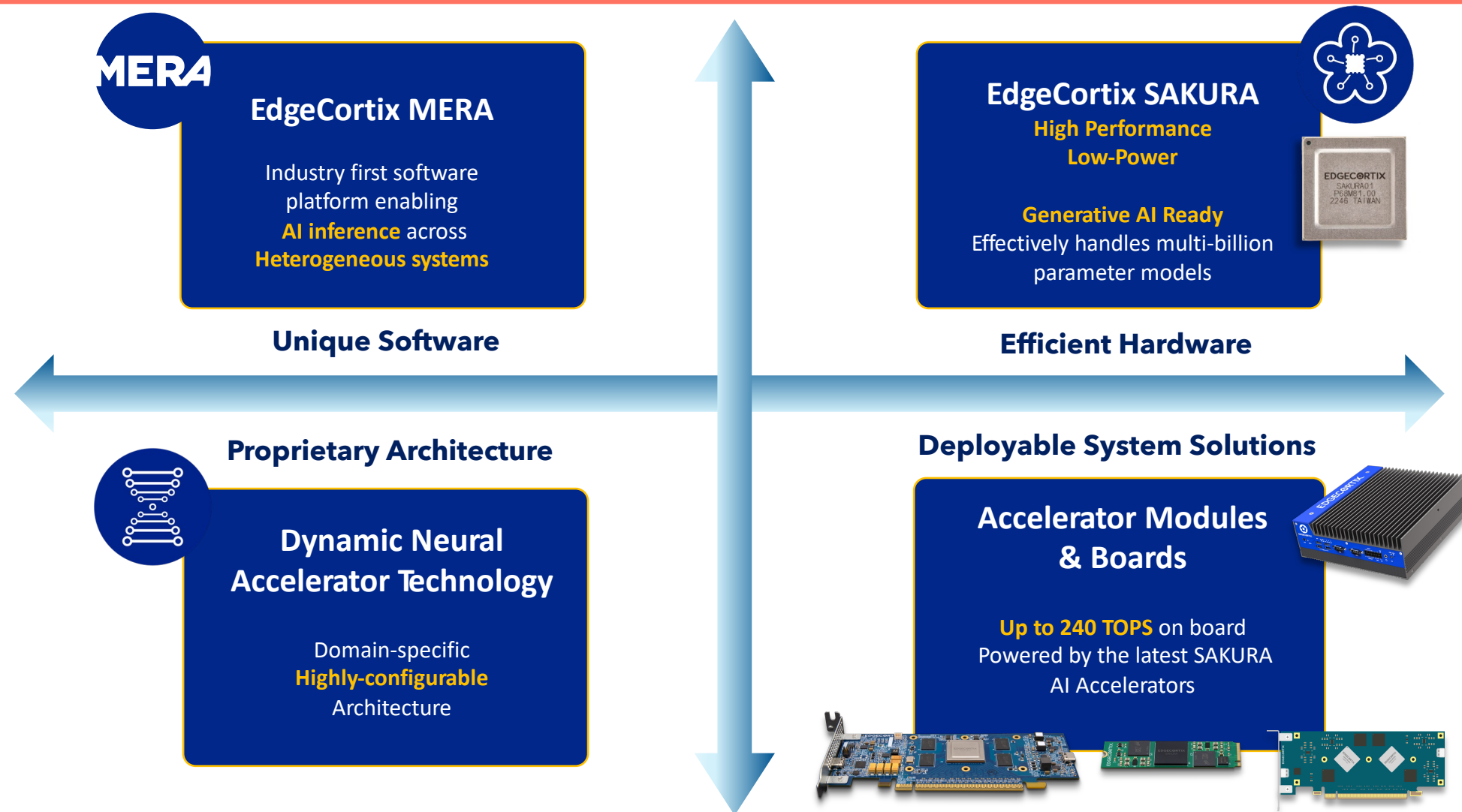


# Breakthrough Efficiency with AI Domain Specific Accelerators



# Software Driven Unified Platform Delivering Highest Efficiency

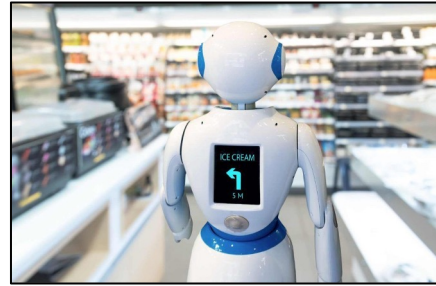
Combining our AI Accelerator with Flexible Software to Deploy Power Efficient Solutions



# Enabling Low-power Generative Edge AI Across Markets



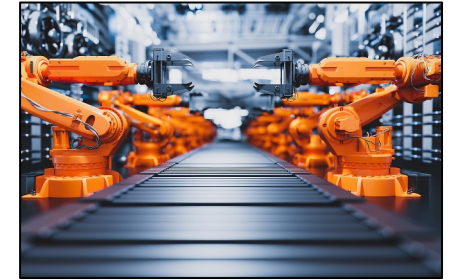
Smart City



Smart Retail



Smart Appliances



Smart Manufacturing



Smart Agriculture



Security



Autonomous Vehicles



Robotics



AI-RAN & Multi-Access Edge Computing (MEC)

## Efficient Edge AI Processing

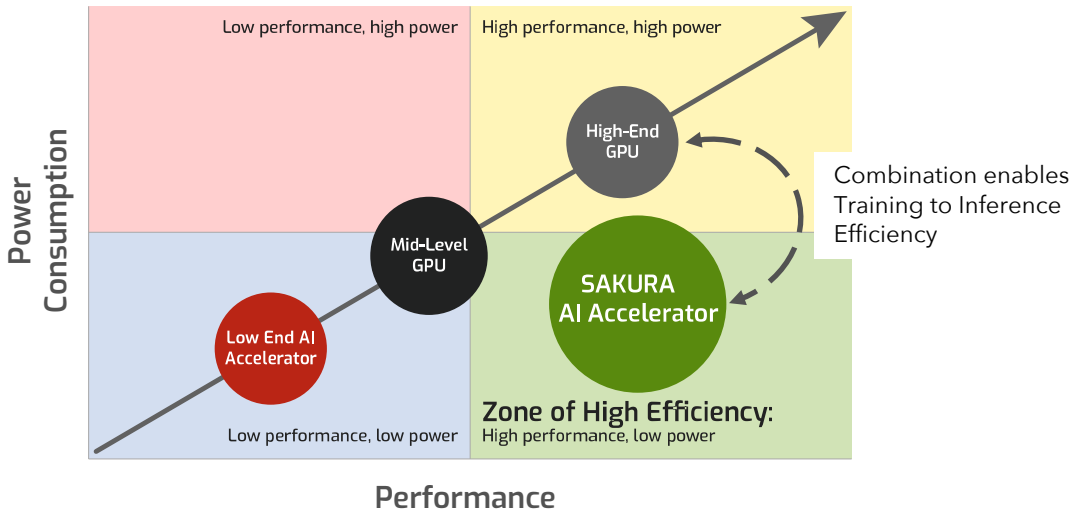
- Natural Language Processing
- Object Recognition
- Person Recognition
- AI enabled RAN
- Segmentation
- Defect Identification
- Obstacle Avoidance
- Signal Processing /MEC



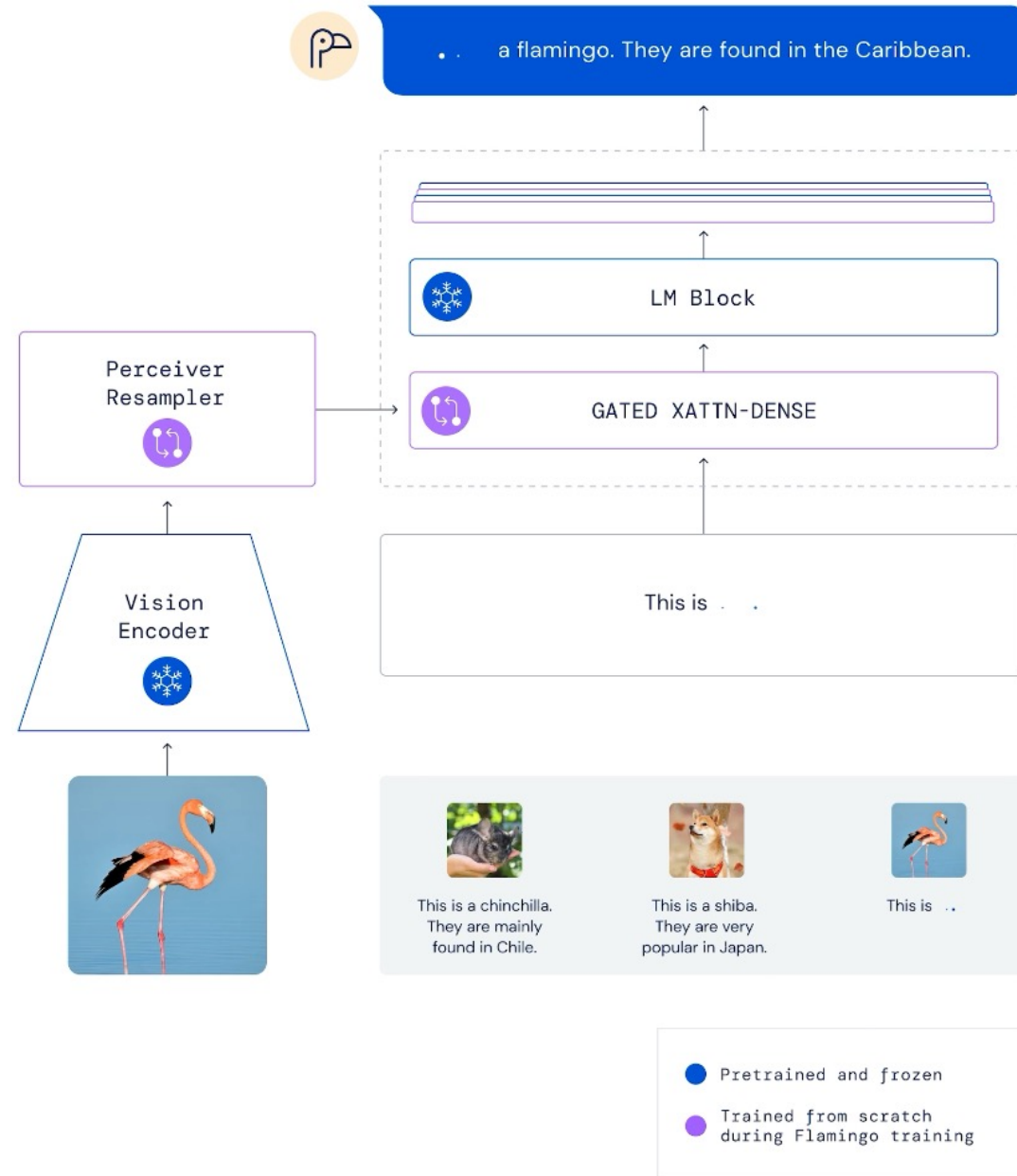
Aerospace & Defense

# With Next-gen SAKURA: Shaping the future of Energy-efficient Generative AI At the Edge

Best-in-class Performance Efficiency



credit: deepmind.google



Delivering trustworthy solutions for the future of AI.

# EDGECORTIX<sup>TM</sup>

Pioneering the Future of Connected Intelligent Edge