**Redefining What's Possible in Edge AI with the ECS-DoT**

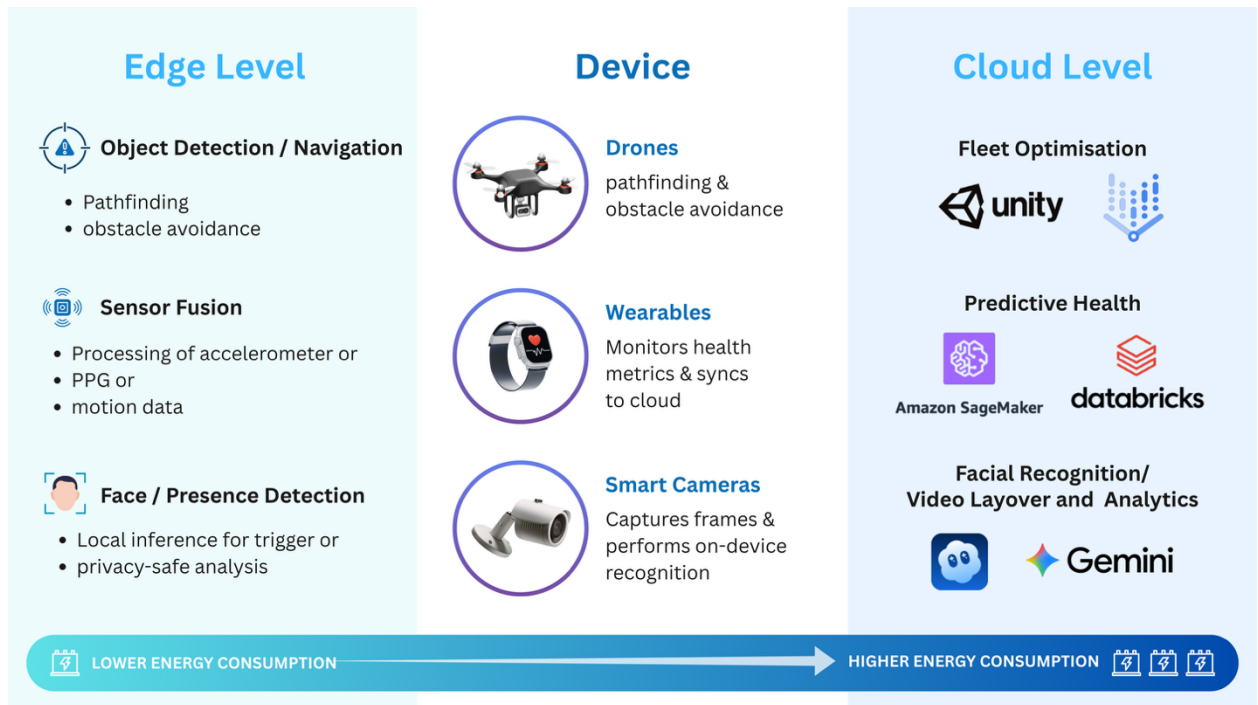**AI Model – the Shift from Cloud-Based to the Micro Edge AI**

We are living in what many refer to as an artificial intelligence (AI) Summer – a golden era marked by rapid innovation, soaring investments, and the seamless integration of AI into daily life. This transformative wave was set into motion in 2012, when deep learning achieved a breakthrough in the ImageNet competition, dramatically outperforming traditional approaches. That moment ignited an unstoppable chain reaction: AlphaGo's triumph in 2016, the rise of transformer architectures in 2017, and the global mainstreaming of generative AI in late 2022, when ChatGPT alone reached over 100 million users in just two months.

Since then, AI has permeated almost every digital interaction. It curates what we see on TikTok and Instagram, personalizes our shopping on Amazon, helps us navigate cities through Google Maps, powers real-time translation, enables biometric authentication, and flags fraud in financial systems. AI has shifted from being a lab-bound novelty to a daily utility, silently shaping our choices, habits, and behaviors. For researchers, it enables breakthroughs across physics, biology, and climate science. And for everyday users, it's a quiet companion embedded in everyday tasks.

This massive growth, however, has largely been powered by centralized cloud AI – systems that rely on sending data from edge devices (phones, sensors, wearables, cameras, vehicles) to distant servers for processing, and then returning the results. While this model has served us well, it carries increasing costs: latency, energy consumption, bandwidth limitations, privacy concerns, and dependence on constant connectivity.

So what's next?

The next logical leap is to bring AI closer to the data source itself—to empower devices to think where the data is born. This is the vision of edge AI – a decentralized, efficient, and context-aware form of intelligence that lives on the device, enabling fast, private, low-power decision-making – with minimal dependence on the cloud.

And yet, for all the excitement, edge AI hasn't fully arrived. Despite remarkable strides in ultra-efficient chips, model compression, and embedded toolchains, most edge deployments remain limited to simple use cases – wake-word detection, object presence, step counting. The leap to truly intelligent, adaptive, multimodal edge devices – capable of running sophisticated AI models while consuming milliwatts or less – remains out of reach for most platforms.

This is where EMASS enters the story. At EMASS, we believe that solving Edge AI's toughest challenges requires a radical rethink of the entire compute stack – from the physics of the transistor to the logic of the application. This is the essence of our "Atoms-to-Apps" philosophy:

- We start with the requirements of real-world applications,

- We co-design hardware and software to work together seamlessly; and,

- We push the frontier of device physics – integrating emerging nanotechnologies, while creating architectures that address the practical limitations of both new materials and conventional approaches.

Our first embodiment of this vision is the ECS-DoT chip – a lean SoC optimized to process data only when needed, dramatically extending what's possible under milliwatt-level power budgets. With it, we aim to prove that true edge intelligence doesn't start in the cloud, it starts with semiconductor physics and is resolved when people use their edge AI products – it goes from atom to app.
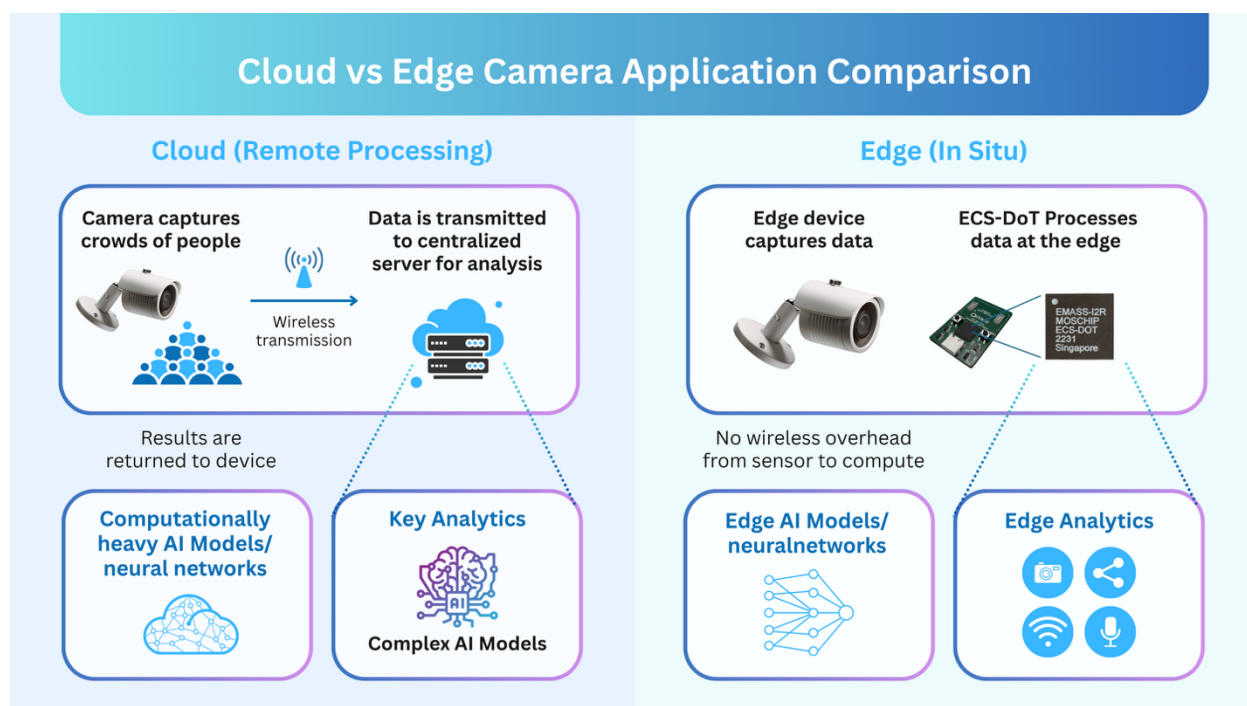
**Edge AI is Fundamentally Different**

AI, as we know it today, is resident largely in the cloud. In the conventional model, data is captured by edge devices (like cameras, phones, or wearables), transmitted to remote data centers where AI models process it, and then the results are sent back to the device. This works well when speed isn't critical and connectivity is stable. But the cloud has limits.

Edge AI fundamentally differs from cloud-based AI by processing data directly on the device rather than sending it to remote data centers. Let's consider a vision-based AI application – say, a real-time security camera tasked with person detection. A typical cloud-based solution captures high-resolution images (e.g., 1080p) and transmits them (a stream that might require megabytes per second of bandwidth) to cloud servers (e.g., AWS GPUs running YOLOv5), achieving high accuracy (~92%) but at the cost of latency (150–300 ms), high energy consumption (~2–3 joules per inference), and potential privacy risks, especially in regulated environments like healthcare or smart homes.

Edge AI addresses these challenges by deploying significantly smaller and highly optimized AI models (e.g., quantized MobileNet-SSD) directly onto specialized low-power chips like EMASS's ECS-DoT or Google's Edge tensor processing unit (TPU). Although these models trade off some accuracy (82–88%), they drastically reduce latency (<10 ms), reduce energy consumption to a few millijoules per inference, and eliminate the need for data transfer entirely, achieving real-time, energy-efficient, and privacy-preserving inference.

Advances in model compression, quantization, and hardware-software co-design have made this approach practical, unlocking entirely new classes of applications previously constrained by power, latency, or privacy.

Edge AI is not just a matter of where the computation happens. It's a fundamental re-architecture of the AI stack with smaller, specialized models; architectures optimized for energy, memory, and latency; and co-designed hardware that can operate within tight constraints, often without active cooling or constant power.

**Cloud vs Edge Camera Application Comparison**

**Edge Is a Spectrum: From Wall-Plugged to Microwatt-Scale Intelligence**

Edge AI exists because the real world doesn't tolerate the delays, costs, and privacy risks of streaming every bit of sensor data off to a distant cloud. When milliseconds matter – when a car is braking, on a factory inspection line, or in a critical health monitor – you can't afford network jitter or expensive bandwidth. Regulations and customer trust demand that personal video and health data stay on-device. And since the infrastructure for capturing that data (cameras, drones, wearables) already exists, it makes economic and logistical sense to "piggyback" AI right where the data is born.

However, edge AI is not one-size-fits-all. Different applications impose vastly different constraints on power, compute, memory, latency, and connectivity. A warehouse gateway can draw tens of watts and host gigabytes of DRAM, while a wearable sensor may have only a few millijoules of energy per inference and a few hundred kilobytes of memory. To address this spectrum, we categorize edge AI into Macro, Meso, and Micro tiers, each calibrated to the unique trade-offs and opportunities at its end of the continuum.
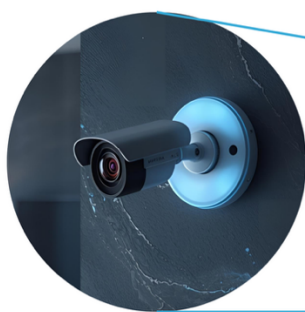
At the macro edge, AI-powered gateways and mini-data-centers sit on factory floors or in server closets. Plugged into the mains and outfitted with 50–300 TOPS (Tera Operations Per Second) of INT8 acceleration and gigabytes of memory, they run large neural networks like YOLO v7 for video analytics or BERT for on-site natural-language processing. For example, Toyota's adoption of an NVIDIA Jetson AGX Orin cluster on its assembly lines reduced defect-inspection delays by 30%, all while keeping sensitive production line data behind the factory firewall and maintaining sub-10 ms inference times.

Moving down the power curve, the meso edge lives inside the devices we carry or fly. Smartphones, AR headsets, and inspection drones—powered by NPUs such as Qualcomm's Hexagon or Apple's Neural Engine, which deliver 5–30 TOPS of compute within a 0.2–5 W envelope. These platforms run compact models like MobileNet-SSD for face unlock, EfficientNet-Lite for real-time translation, or tiny Transformer variants for contextual assistance. Each inference completes in 5–20 ms and consumes only a few millijoules, letting drones survey warehouses, headsets translate conversations, or phones process images without draining the battery.

At the far end sits the Micro Edge: always-on sensors, smart tags, and wearables that must last months or even years on coin cells or harvested energy. Here, chips like EMASS's ECS-DoT squeeze 0.1–1 TOPS and a few hundred kilobytes of SRAM into a sub-10 mW package. They execute TinyML workloads such as wake-word detection, presence sensing, vital-sign anomaly spotting, in under 10 ms and at just 1–10 µJ per inference. This whisper-quiet, ultra-efficient intelligence enables truly autonomous sensing in environments where changing batteries is impractical or impossible.

Together, macro, meso, and micro edge form a seamless continuum – each tier tuned to balance latency, energy, cost, and privacy in its own way. In the pages that follow, we'll zoom in on the micro edge, where EMASS's Atoms-to-Apps co-design philosophy and the ECS-DoT architecture break through previous limits, unlocking a new era of always-on, ubiquitous intelligence.



**Battery Power Landscape Across Edge Devices**

**Macro**
Non-battery concious devices
robotics, surveillance cameras

**Meso**
Medium-power devices
smartphones, drones, gateways

**Micro**
Ultra-low-power sensors
wearables, hearables

**EMASS Holistic Approach**

Atoms-to-Apps embodies EMASS's end-to-end approach to systems design. The semiconductor industry has driven integration down to sub-micron dimensions, a level where features can in fact measure mere atoms across or deep. Atoms-to-Apps starts at the device level and flows through circuit design, micro-architecture, compilers, runtimes, and application optimization. If the objective is ultra-low-power, always-on inference with high fidelity, devices, chips, models and applications cannot be treated as independent silos.

Orchestrating design decisions across conceptual layers – and propagating constraints and metrics both up and down the stack – unlocks efficiency and capability gains far greater than optimizing any single layer in isolation.

Below we highlight the main pillars that EMASS' approach is built on.

**Nanodevices — Concise System Implications**

Emerging nanodevices are the physical levers that change what is possible at the micro edge. For micro edge AI the two most relevant are storage devices:  resistive memories (RRAM / ReRAM) and magneto/ferroelectric devices (MRAM / FeFET). Each class of memory trades density, persistence, read/write energy, and endurance differently – and those tradeoffs directly determine circuit choices, array organization, and software strategies.

RRAM / ReRAM (scalable, BEOL-friendly, MLC capable): RRAM is attractive for micro-edge designs because it scales well in area, can be integrated into the back end-of-line (BEOL) of standard CMOS processes, and supports multi-level (multi-bit) storage per cell. Those properties make RRAM a natural candidate for neural networks (especially those with very dense on-chip weight banks) and for analog Compute-in-Memory (CIM) implementations that similarly work well in neural networks  (by reducing the need for energy-intensive transfers of data off and on chip). In practice, the array-level picture is complex – device variability, nonlinear I–V characteristics, ADC and peripheral overhead, and endurance/retention tradeoffs must be managed. When peripheral, calibration and error-management costs are accounted for, RRAM can deliver large wins for inference-dominated workloads that benefit from persistent, dense weights.

MRAM / FeFET (robust persistence, low-energy updates): MRAM offers fast reads, strong retention and high endurance, making it ideal for persistent caches, control state, and small weight buffers. FeFETs complement this role by offering low-energy writes and CMOS-friendly voltages, enabling near-instant resume with minimal cold-start penalty. Both are preferred where frequent small updates or high read rates are expected and where endurance budgets can be managed.

**Logic Devices and Integration**

Advanced logic structures, such as Gate-All-Around (GAA) transistors and  Carbon Nanotube Field - Effect Transistors (CNFETs) can lower leakage and push near-threshold

energy efficiency, while heterogeneous integration structures and techniques (such as interposers and 3D stacking, respectively) can dramatically shorten interconnects and increase bandwidth-per-watt. Both yield system advantages ,but bring supply-chain, thermal and BOM costs that must be justified by end-to-end gains.

### Array and Peripheral Realities

Device numbers alone are misleading – ADC resolution and count, sense drivers, calibration, Error Correction Codes (ECC) and periphery area often dominate energy and area in CIM designs. Any decision to adopt RRAM/CIM must amortize peripheral costs across sufficient on-chip density or workload patterns (e.g., large, read-dominant weight banks).

### Practical System Implications

Edge AI is an inherently Atoms-to-Apps endeavor because decisions made at either end will almost certainly have implications for the other. The results desired from any application will dictate the best AI techniques to use, which in turn will dictate implementations at the circuit and device level. Here are just some examples of Atoms-to-Apps design considerations:

- Match device choice to workload: inference-dominated, read-heavy applications are the best fit for persistent RRAM weight storage; frequently updating workloads favor MRAM/FeFET/SRAM (depending on macro sizes and offerings by foundries).

- Use a hybrid memory hierarchy:

    o   SRAM for hot working sets

    o   MRAM/FeFET for state and small buffers (if permissible by foundries)

    o   RRAM for dense, persistent weight banks when peripheral costs are amortized.

- Exploit RRAM's BEOL compatibility and MLC capability to pack large models or multiple quantization levels on-chip, but design compilers and runtimes to tolerate variability and non-idealities.

- Quantify non-ideal costs (ADC energy, calibration, ECC, endurance management, packaging) early in the design flow.

ECS-DoT adopts a pragmatic hybrid approach: leverage RRAM's density and BEOL integrability for persistent weight storage where justified, use MRAM/FeFET for fast state and small caches (if needed and adequate macro sizes are offered by foundries), and retain SRAM for critical hot paths. The design rule is simple – commit to RRAM/CIM only when a full stack budget (including ADC/peripheral and endurance costs) shows clear, repeatable system advantage.

**Managing Power is Critical**

The central circuit principle for ultra-efficient micro-edge AI is straightforward and practical – only the circuits that must be doing work should consume meaningful energy – devote energy only to what must be active.

In sparse-duty, always-on systems, the system spends most of its life idle. This means the single largest lever for efficiency is to shrink idle power to a vanishing fraction of active power. When done correctly, this approach converts short bursts of inference into a cost that is negligible compared with long idle intervals – enabling months or years of battery life on tiny cells, which should open use cases that were previously impossible.

Putting that principle into engineering practice yields three immediate and interdependent priorities that must drive design choices across device, circuit and architecture.

First, minimize the always-on budget before anything else. For micro-edge nodes, idle leakage typically dominates total energy for sparse workloads. Hence minimizing or even eliminating idle power altogether is critical to achieve significant power savings. Aggressive power minimization through gating and voltage scaling are a necessity in this context.

Second, make wake costs predictable and far smaller than saved idle energy. Any gating or persistence strategy only helps when the energy to resume is meaningfully lower than the idle energy it replaces. Design to a clear break-even point for your expected duty cycle, and favor small, always-available sensor front ends plus persistent state so resume is fast, low-energy, and repeatable in the field.

Third, trade granularity against complexity with discipline. Fine-grained power domains and retention islands unlock the biggest savings, but each added domain increases control logic, verification scope, and area. Match domain granularity to the application's active set and duty pattern: coarse domains where activity is frequent, fine domains where long-idle intervals dominate.

Finally, measure everything and hold to targets. EMASS treats the above as constraints, not recipes – implementation options (power switches, retention islands, persistent storage) are chosen only when they demonstrably move the system toward agreed, end-to-end targets under realistic manufacturing and cost assumptions. That discipline

<div align="center">principle → target → justified implementation</div>

is what turns the "direct energy" idea into reliable, deployable micro-edge products.

**Microarchitecture**

At micro scale the microarchitecture is not about raw peak TOPS — it is about how data moves. AI ordinarily entails frequent transfers of data between logic and memory. At the Macro level, this is inconvenient, but as you move through Meso to Micro, the amount of energy consumed by these operations becomes problematic The design objective, then, is to exploit locality and reduce data footprint so that energy spent moving bytes is minimized

and energy is spent doing arithmetic. To that end, EMASS adopts a memory-centric microarchitecture where compute and memory are interwoven.

The AI acceleration module we built is very efficient; it is a compact, optimized acceleration fabric purpose-built for micro-edge workloads. This fabric is further enriched by a handful of specialized blocks for non-linear ops and accelerated support for algorithmic optimizations performed at the application level, and near-memory interface, with scalable interconnects to on-chip memory banks.

EMASS' microarchitecture module also integrates a lightweight device-management IP that addresses practical nanodevice challenges with minimal runtime cost. This IP provides low-overhead services such as per-bank calibration offsets, ECC/sensitivity compensation, wear-aware allocation and write-minimization policies, and compact health telemetry. Critically, these functions run locally, and are sized to avoid negating the energy gains from near-memory or CIM approaches. In short: device non-idealities are managed as part of the architecture, not as an afterthought.

**Application Layer**

At the application layer the problem is simple to state and fiendishly hard to solve in practice: deliver the required task-level accuracy and latency while fitting the model and its working set into extremely tight memory and energy budgets.

Success requires treating algorithms as first-order design levers, not as artifacts to be squashed into hardware after the fact. That means selecting and combining model transforms that reduce memory footprint and active computation without introducing extra memory traffic or unnecessary operations that would defeat the energy budget.

Common, widely-used techniques – described here as engineering primitives rather than proprietary recipes – include quantization and mixed precision (per-channel / per-layer where useful); structured pruning and block sparsity that allow whole lanes or banks to be gated; operator fusion and reordering to eliminate intermediates; tiling and streaming so only a small working tile resides in SRAM; knowledge distillation to produce compact students; early-exit cascades so heavy stages run only when needed; low-rank / separable factorizations to reduce memory and compute; and on-demand compression schemes that trade a small decompression cost for dramatically reduced persistent-store transfers. Each technique is chosen not merely to save ops but to reduce bytes moved and to create large, hardware-friendly opportunities for power gating.

Crucially, these algorithmic choices are not independent of the microarchitecture or the device layer; they are co-designed. The acceleration module exposes primitives (bitwidth controls, sparsity masks, fusion hooks, bankable memory tiles) so the compiler maps transforms onto the hardware with minimal control overhead. Simultaneously, the compiler schedules memory access patterns to maximize locality, avoid frequent writes to

endurance-sensitive NVM, and amortize ADC/peripheral costs when analog CIM is used. In short: algorithms become hardware-aware (they produce locality, sparsity and predictable access), the architecture is algorithm-aware (it supports those primitives efficiently), and the runtime enforces device-aware policies (wear-aware allocation, amortized maintenance, on-demand decompression) so non-idealities remain visible and manageable. In short: algorithms are written for the device; the architecture provides low-overhead primitives to implement them; and the runtime closes the loop to deliver predictable energy, latency and lifetime.

**Introducing the ECS-DoT – A New Architecture for Edge Intelligence**

The ECS-DoT is the first embodiment of EMASS's Atoms-to-Apps philosophy – an SoC designed from the ground up for ultra-low-power, always-on intelligence at the micro edge. Unlike conventional SoCs that bolt on AI accelerators, ECS-DoT integrates compute, memory, and sensor interfaces into a single architecture optimized for real-time, local machine learning.

By fusing architectural efficiency with practical deployment features, ECS-DoT extends the principles discussed throughout this paper – memory-centric design, hybrid device integration, and power-aware microarchitecture – into a product that can be designed into next-generation devices today.

Architecture highlights:

- Power efficiency – operates at 0.1–5 mW per inference, enabling continuous, always-on sensing in battery-constrained devices.;

- Latency / real-time responsiveness – achieves <10 ms per inference, delivering instant responsiveness without cloud or host processing;

- Energy per inference – consumes just 1–10 µJ per inference, allowing complex AI processing locally without draining batteries;

- Multimodal sensor fusion – handles inertial sensors, microphone, and camera inputs in parallel at milliwatt-scale for richer contextual awareness;

- Integrated memory and on-device processing – includes up to 4 MB on-chip memory for fully on-device model execution with no off-chip DRAM.

While many companies have introduced NPUs and microcontrollers that claim "edge AI" capabilities, most remain constrained by high power draw, limited memory, and sluggish responsiveness. The ECS-DoT stands apart by delivering significantly lower latency (<10 ms) and dramatically lower energy per inference (1–10 µJ)—making true always-on intelligence feasible in battery-powered devices. The table below highlights how ECS-DoT

compares to leading alternatives, underscoring its unique position in enabling real-time, multimodal, and sustainable AI at the micro edge:

| Feature / Metric | EMASS ECS-DoT | Competitor A | Competitor B |
|---|---|---|---|
| Power per Inference | 0.1-5 mW | 5-100 mW | 30-150 mW |
| Latency | <10 ms | 10-15 ms | 10-100 ms |
| Energy per Inference | 1-10 µJ | 30-150 µJ | 100-2,000 µJ |
| On-Device Memory | Up to 2MB SRAM + 2MB MRAM/RRAM | up to 1MB SRAM | 2MB SRAM with optional 2MB MRAM |
| Multimodal Sensor Fusion | Yes, milliwatt-scale | Limited | Limited |
| Always-On Viable? | Yes | No | No |

**Application Enablement – Where ECS-DoT Creates Value**

ECS-DoT's architecture directly translates into system-level benefits that open new opportunities across industries:

- Wearables – Extends smartwatch and fitness tracker battery life by days through efficient biosignal, gesture, and activity inference. Medical wearables can run 24/7 anomaly detection without frequent recharging, improving compliance and safety.

- Drones – Every milliwatt saved translates into longer flight time and payload capacity. ECS-DoT enables real-time navigation, obstacle avoidance, and predictive maintenance onboard, independent of cloud connectivity.

- Industrial IoT & Robotics – Delivers predictive maintenance by analyzing vibration and acoustic signatures at the edge, reducing downtime and service costs. Millisecond-level inference supports faster robotic control loops and safer factory automation.

- Smart Devices – Cameras, remotes, and security systems process data locally, reducing bandwidth use and preserving privacy. Devices wake intelligently only when needed, improving responsiveness while lowering energy consumption.

- Healthcare & MedTech – Portable diagnostics and patient monitoring systems gain continuous operation under tight power budgets, enabling real-time vital sign monitoring and anomaly alerts in hospital and home settings.

- AgTech & Environmental Sensors – Ultra-low power enables season-long deployments on coin cells or solar harvesting. ECS-DoT processes soil, water, and air quality data locally, reducing wireless transmission and enabling autonomous field sensing.

## Conclusion – from Architecture to Applications

Edge AI's future depends not just on smaller models or lower power chips but on architectures built intentionally for the micro edge. ECS-DoT is that architecture—extending Atoms-to-Apps design principles into a chip that delivers real benefits in the field: longer drone flight times, longer-lasting wearables, safer factories, smarter healthcare, and sustainable environmental monitoring.

By collapsing complexity into one ultra-efficient package, ECS-DoT not only redefines what's possible today but also signals the direction of the industry: always-on, everywhere, and fundamentally local intelligence.