

# Device First Continuum AI (DFC-AI): Realizing Human-Like AI

Siavash M. Alamouti, Fay Arjomandi, Michel Burger

mimik Technology, USA

**Abstract.** This study introduces Device First Continuum AI (DFC-AI), a transformative architecture within the Hybrid Edge Cloud paradigm designed to address the limitations of traditional cloud-centric artificial intelligence across diverse applications. DFC-AI prioritizes the deployment of intelligent agents, built on a microservices framework, that originates and primarily resides on end devices, extending to gateways and cloud servers as needed. This Device-First approach is essential for enabling real-time decision-making and personalized experiences for both industrial and consumer applications, particularly in scenarios demanding low latency, operation in disconnected environments, and efficient management of massive data streams. The study highlights the fundamental challenges of relying solely on centralized cloud or basic edge computing models, including prohibitive bandwidth costs, energy inefficiency, and compromised user privacy. By embedding intelligence at the device level, DFC-AI overcomes these limitations, fostering autonomous operation, seamless collaboration among devices, and substantial reductions in operational overhead, moving us closer to realizing the potential of truly human-like artificial intelligence in machines. Through illustrative examples spanning various sectors, this study demonstrates the potential of DFC-AI to unlock a new era of holistic, responsive, and user-centric intelligent systems, paving the way for innovative applications and enhanced digital experiences in an increasingly connected world.

**Keywords:** Device First Continuum AI, Hybrid edge cloud, Cloud servers, Gateways, Artificial intelligence

## 1 Introduction: Beyond Siloed Machines, Towards Holistic Intelligence

Current approaches to Artificial Intelligence are hitting a wall. Siloed, cloud-centric models and even basic "edge computing" architectures [1], [2], [3] struggle to deliver truly holistic, efficient, and responsive intelligence, particularly in the real world, where network connectivity is unreliable, and context is paramount.

The prevailing cloud-centric paradigm [4] in Artificial Intelligence (AI), while historically dominant and mirroring early centralized computing models, is now proving insufficient to deliver the truly holistic and real-world intelligence required for future applications. Moving beyond these siloed models, a new architecture is essential:

the Device-First Hybrid Edge Cloud [5], [6], inspired by the human brain's distributed and efficient intelligence [7]. This Device-First Hybrid Edge Cloud architecture, mirroring the human brain's interconnectedness, prioritizes integrated, responsive, and collaborative intelligent systems.

Consider the limitations of most AI systems today. Robots confined to repetitive tasks in factories, smart assistants isolated within individual devices, or even self-driving cars struggling with unexpected situations. These exemplify siloed AI. While possessing narrow capabilities, they often lack crucial aspects of human-like intelligence: the ability to grasp context, adapt dynamically to changing environments, learn on the fly, and collaborate effectively. They represent fragmented intelligence, not the holistic, integrated intelligence we need. This fragmented approach to AI is increasingly inadequate for realizing its full potential.

To unlock a world where AI truly empowers us and solves complex real-world problems, we must move beyond these limitations and embrace Holistic Machine Intelligence. This means developing AI systems that emulate the key attributes of human intelligence: being comprehensive, adaptable, and inherently collaborative. Envision intelligent ecosystems, not just isolated smart devices. Imagine factories where robots collaborate and learn dynamically, drone fleets orchestrating deliveries in real-time, and smart homes that truly understand and respond to human needs.

This is not a vision for the future but absolutely plausible today, and the right starting point for our future evolution. Failing to embrace this paradigm shift now risks significantly delaying the operationalization and widespread benefits of truly intelligent, AI-enabled systems.

To achieve this Holistic Machine Intelligence, we must go beyond siloed approaches. Using the outdated cloud models for the new agentic world is like using square wheels to build a cart. We need to rethink AI architecture, moving from cloud-heavy, centralized models to something distributed, device and offline-first, and truly integrated. Just like our brains use billions of connected neurons to be smart, Holistic Machine Intelligence must start from end devices working together across a flexible continuum, all the way to cloud servers.

This study will explore why the old models are not economically and operationally feasible, and why we urgently need to leave behind these AI silos and take the path to Realizing Human-Like AI, inspired by the example of human-like comprehensive intelligence.

## **2 The Human Blueprint: Holistic Intelligence in Biological Systems**

To truly understand the potential of Holistic Machine Intelligence, we can look into the most compelling example of intelligence we know: ourselves. The human brain [8], and indeed the entire human nervous system, is a marvel of holistic design. It's not a collection of isolated modules. It's a profoundly interconnected and collaborative system. Consider how effortlessly we integrate vast amounts of information from diverse sources: sight, sound, touch, smell, proprioception, to form a unified, coherent understanding of our environment. We don't experience the world as fragmented data streams. We perceive a seamless, integrated reality.

Furthermore, human intelligence is characterized by its remarkable adaptability and efficiency. We can learn new skills, adjust to changing circumstances, and solve novel problems with impressive flexibility and resourcefulness [9]. Our brains achieve complex computations with astonishing power efficiency, far surpassing even the most advanced AI in energy-per-operation. And crucially, human intelligence is inherently social and collaborative. We thrive in teams, share knowledge, learn from each other, and build collective intelligence that exceeds the capabilities of any individual.

This holistic nature of human intelligence includes integration, adaptability, efficiency, and collaboration. It isn't accidental. It's the result of over six million years of evolution and is fundamental to our success as a species. It's the blueprint we should aspire to when designing truly intelligent machines. Current AI, with its siloed architectures and fragmented capabilities, is a far cry from this holistic ideal. But by understanding the principles of biological intelligence, particularly the human brain's distributed and collaborative nature, we can chart a course towards a new generation of AI with Holistic Machine Intelligence that truly mirrors the comprehensive capabilities of the human mind.

### **3 Today's World: Fragmented versus Holistic Approaches**

Look around at the AI already in our world. Smart speakers answer questions, robots vacuum our floors, and apps offer personalized recommendations. These are all impressive feats of technology, yet they largely represent fragmented AI in action. Consider a smart speaker: it's good at voice recognition and accessing cloud-based information, but it operates in isolation. It's unaware of other smart devices in your home, unable to collaborate with your smart thermostat to optimize energy use, or coordinate with your smart lighting system to create a more immersive experience. It's a smart island, not part of a holistic intelligent home.

Think about today's delivery drones. They autonomously navigate routes, but each drone typically operates independently, with limited real-time coordination with other drones or ground traffic systems. Imagine the potential if these drones could communicate dynamically, share airspace information, and reroute collaboratively to optimize delivery efficiency across an entire city. That would be a move towards holistic drone fleet intelligence, far beyond today's siloed operations. Even advanced factory robots, while increasingly sophisticated, often work within fixed, pre-programmed routines, lacking the adaptability to dynamically reconfigure workflows in real-time collaboration with other robots and human workers on the factory floor.

These examples highlight the limitations of fragmented, siloed AI in real-world applications. While cloud-centric and gateway-mediated architectures have enabled significant progress, they inherently encourage this fragmentation. DFC-AI offers a fundamentally different path: one that fosters holistic intelligence by empowering end devices to become active, collaborative participants in an integrated AI ecosystem. By bringing intelligence to the device level and enabling seamless device-to-device and device-to-cloud workflow collaboration, we can pave the way for AI systems that are not just collections of smart fragments but truly holistic, integrated, and collaborative intelligences, capable of solving real-world problems with a level of sophistication that mirrors human ingenuity.

#### **4 Addressing the "Edge Computing" Confusion: Device-First is Not Just About Location or More Expensive Hardware**

You might hear "edge computing" and think Device-First is just another name for it. Let's be clear: DFC-AI is not simply about moving computation to the network edge to save latency or bandwidth, nor is it about adding GPUs to smaller devices like cameras and sensors. These are improvements, but still often mean centralized control, or a siloed paradigm, and miss the point of truly holistic machine intelligence. Even with these steps, we can still get islands of intelligence, not a real collaborative ecosystem.

DFC-AI is about going beyond the location of computing. It's a fundamental re-architecture of AI to consider the entire AI continuum, starting with end devices from the ground up. It's about making individual devices such as smartphones, robots, drones, vehicles, and even sensors to be active, intelligent agents. They get to make some decisions independently, process data locally, and seamlessly collaborate with other devices, gateways and servers in the cloud when needed. We must build a truly distributed, dynamic AI ecosystem where intelligence is built-in at every level, starting right at the end device.

Think of "edge computing" as setting up outposts of a command center or just a bunch of separate smart devices. Instead, we need a network of capable, autonomous entities that can work together without being micromanaged. And just like in human groups, where children, adults, athletes, and thinkers all contribute differently but are equally valuable, we need to create a system where devices of all kinds, from simple sensors to powerful robots, can collaborate and add value, regardless of their hardware. The point isn't just where processing happens, or even device cost and power. It's about how intelligence is designed, distributed, and made to work collaboratively across the whole ecosystem. DFC-AI enables this deeper change, going beyond fragmented "edge" ideas to a truly holistic, collaborative intelligent fabric.

#### **5 Offline-First and Context-Aware Workloads: Why End Devices Need To Be Independent and Smart About Their Surroundings?**

We've established that we must build collaborative, ecosystem-level intelligence, going beyond simple "edge computing". But to truly realize Holistic Machine Intelligence, intelligent systems architectures need two more essential capabilities: Real-time Offline-First operation and Context-Awareness. These aren't just optional features; they are fundamental enablers that unlock the full potential and make them truly human-like.

Think about human intelligence. We don't stop being intelligent when our internet connection drops. We can still navigate, make decisions, and even collaborate in real-time, locally, and even when offline [10]. Imagine if an autonomous vehicle lost all its intelligence the moment it went through a tunnel or into a dead zone! That would be awful. Similarly, in such intelligent systems, end devices need to be Offline-First, meaning capable of functioning intelligently and continuously even when network

connectivity is intermittent, unreliable, or completely absent. This isn't just about convenience. It's about resilience and reliability. Critical applications like autonomous vehicles, factory robots, and emergency response systems must be able to operate offline to ensure continuous, dependable performance in any situation.

Beyond just staying functional offline, human intelligence is also deeply context-aware. We don't operate in a vacuum; we constantly sense and respond to our surroundings. We understand context (location, time of day, social cues, our physical and psychological state, environmental conditions) and adjust our behavior accordingly. For AI to be truly human-like, devices need to be context-aware too. The DFC-AI architecture must excel at this. By direct access to information from the local environment and operating systems, devices should directly sense and interpret their surroundings, understanding their hardware and software resources, physical location, user activity, nearby devices, and real-world conditions, all in real-time. This context-awareness is crucial for effective collaboration.

Imagine robots on a factory floor needing to dynamically adjust their tasks based on the location of other robots and workers, or a smart home adapting lighting and temperature based on the time of day and occupancy of different rooms. This level of dynamic, context-driven behavior is only truly possible with end devices enabled by human-like intelligence. The alternative of using cloud servers or gateways would require transmitting terabytes of data on a daily basis for complex devices such as autonomous vehicles or hundreds of gigabytes by robots and drones.

Offline-First and Context-Aware workloads are therefore essential enablers of Holistic Machine Intelligence. They allow systems to be resilient, reliable, adaptable, and truly intelligent in dynamic, real-world environments, mirroring the key characteristics of human intelligence and moving us closer to realizing the full potential of human-like AI.

## 6 Choices of Architecture for Intelligent Systems

To truly understand the distinct advantages of DFC-AI, it's helpful to compare it directly with other common architectural approaches for intelligent systems. Table 1 provides a concise comparison of three key architectures: Cloud-Based [4], Gateway-Mediated [11], [12], and DFC-AI built using a Hybrid Edge Cloud [5] architecture across a range of requirements essential for realizing Holistic Machine Intelligence. This comparison illustrates the unique strengths of the DFC-AI approach in building truly holistic and human-like AI systems.

**Table 1.** Comparison of cloud native architectures for intelligent systems

Attribute	Cloud-Based AI	Gateway-Mediated AI	DFC-AI
<b>Collaboration &amp; Interoperability</b>	Low (Cloud-Mediated)	Moderate (Gateway Domain)	<b>High (Peer-to-Peer)</b>
<b>Offline Operation</b>	Poor (Cloud Dependent)	Limited (Gateway Cache)	<b>Excellent (Device-Centric)</b>
<b>Contextual Awareness</b>	Poor (Cloud View)	Poor (Gateway View)	<b>Excellent (Device View)</b>

<b>Energy Efficiency</b>	Lower Efficiency (Cloud datacenters consume high power)	Medium Efficiency (Gateways improve regional efficiency)	<b>Highest Integrated Efficiency (80% on average compared to cloud)</b>
<b>Device Heterogeneity Support</b>	Challenging (Uniformity)	Moderate (Domain Abstraction)	<b>Excellent (Native Support)</b>
<b>Cloud-Hosting Costs</b>	Highest Cloud Costs (Cloud servers for most processing)	Potentially Reduces Cloud Hosting Costs	<b>Lowest Cloud Costs (Up to 95% less as Devices handle local tasks)</b>
<b>Operational Effectiveness</b>	Less Effective for Real-Time, Context-Aware Apps	Marginally More Effective	<b>Highly Effective for Real-Time, Context-Aware Apps</b>
<b>Update Complexity</b>	Complex (Centralized)	Simplified (Domain)	<b>Simplified (Distributed)</b>

## 7 Understanding the Architectural Trade-offs: Key Insights from the Comparison Table

This table evaluates cloud-centric, Gateway-Mediated, and DFC-AI architectures not as isolated solutions, but as interconnected components within a holistic Hybrid Edge AI ecosystem. It highlights the distinct strengths and trade-offs of each approach from a system-level perspective, focusing on how each architecture contributes to or hinders the creation of a truly integrated, efficient, and responsive AI ecosystem. Understanding these nuanced differences across deployment focus, agent architecture, integration capabilities, efficiency, and cost is crucial for choosing the optimal architecture, or combination of architectures, to realize the full potential of Hybrid Edge AI and achieve human-like intelligence in real-world applications. The table summarizes these key architectural differences across various critical features.

The comparison shows the distinct trade-offs inherent in different Hybrid Edge AI architectural approaches. Cloud-centric architectures, while offering immense global processing power, sacrifice latency, efficiency, and real-time responsiveness. Gateway-Mediated approaches offer a localized improvement in responsiveness and perhaps cost, but still rely on centralized hubs and can limit overall dynamism. In contrast, DFC-AI is uniquely positioned to deliver on the promise of truly integrated, efficient, and responsive Hybrid Edge AI. By prioritizing intelligence at the device level and enabling seamless continuum-wide deployment, DFC-AI offers superior contextual awareness, power efficiency, and low latency, making it highly effective for real-time, context-aware applications and paving the way for a more sustainable and scalable AI future.

## 8 Energy Efficiency Through Microservices at the End Device Edge: Emulating the Brain's Energy Conservation at the Source of Context

Let's zoom in on energy efficiency, a paramount consideration when we talk about Edge AI, especially at the end device located at the very edge, where power is often limited,

and efficiency is crucial for truly pervasive intelligence. When we emulate human intelligence, we must also emulate its remarkable energy efficiency, particularly at the "edge" of our own bodies: our brains and sensory organs [13]. Microservice architectures [14], [15], and [16], deployed directly on the end device, mirror this elegant energy conservation strategy.

Just as our brains have specialized neural circuits for different tasks, microservice agents are built from smaller, specialized components. And just like our brains, these microservices enable:

- 1) **Selective Activation: "Neural Circuits on Demand":** Think of microservices as specialized neural circuits. Only the circuits actively needed for a particular task are powered up and consuming energy. Microservices handling object recognition might be active when the device is "seeing", while microservices handling audio processing are activated when the device is "hearing". Idle microservices, like dormant neural pathways, can be put into a low-power state, minimizing energy waste. It's like our brain selectively illuminating only the relevant areas for a given thought, instead of firing up the entire cortex for every task.
- 2) **Optimized Resource Allocation: "Right-Sizing Brainpower":** Microservices allow for precise control over resource allocation like our brain allocating just the right amount of "mental horsepower" for each task. If you're just walking, your brain uses a different level of processing than when you're solving a complex problem. Similarly, microservices prevent over-provisioning resources and wasting energy on unused capacity. It's about using exactly the right neural "engine" for each job, big or small, instead of always running a massive, overpowered brain at full throttle.
- 3) **Event-Driven Awakening: "Reflexive Responses":** Microservices can be designed to be event-driven, much like our reflexes. They lie dormant, using minimal power, until a specific sensor reading, user input, or external event acts as a trigger, "waking them up" to perform their specialized task. Imagine AI services only "firing up" when something relevant actually happens in the device's environment, instead of constantly running and consuming energy even when there's no stimulus. Think of how your reflex to pull your hand away from a hot stove is instantaneous and energy-efficient, powered only when needed.

For end devices, especially battery-powered ones (wearables, remote sensors, mobile robots), this microservice-driven power efficiency isn't just about cost savings. It's about mimicking the essential biological efficiency that allows complex intelligence to thrive in resource-constrained environments. It's about enabling truly pervasive, long-lasting, and practically deployable Edge AI.

DFC-AI microservices-based architecture is specifically designed for power efficiency. Edge-native microservices are lightweight and resource-optimized, enabling devices to participate intelligently in the continuum without draining battery life or requiring excessive processing power. This distributed efficiency, combined with strategic use of gateways and cloud resources only when necessary, results in a Hybrid

Edge AI system that approaches the remarkable energy efficiency of biological intelligence, paving the way for truly sustainable and scalable AI deployments.

## 9 Collaborative Compute and Capability Transfer with Microservices: Building a “Social Brain” *Right at the Edge of Context*

Now, let's explore collaboration, envisioning a “social brain” [17] not just at the network periphery, but truly distributed and interconnected at the end device edge itself. This is where the power of microservices truly shines, enabling direct, peer-to-peer collaboration between end devices, the very locations where context is richest, and collaboration is most impactful. Because microservices are modular, independent, and crucially discoverable within a network, they mimic the way neural networks in a human brain can connect and exchange information.

Let's consider a robot warehouse example and flesh it out within the edge collaboration paradigm:

- 1) **The Scenario: The Collaborative Warehouse "Brain" at the End Device Edge:** Imagine a warehouse teeming with robots from diverse manufacturers, each now conceived as an intelligent "node" in a distributed warehouse brain. Robot A (Manufacturer X) is like a highly specialized “visual neuron” excelling at object recognition. Robot B (Manufacturer Y) is the “navigation neuron,” brilliant at path planning. Robot C (Manufacturer Z) is the “communication neuron,” adept at inter-robot communication and task coordination. In a traditional, siloed approach, these robots would operate as isolated units, limited to their own hardwired "neuronal" pathways.
- 2) **The Microservice Revolution: "Neurons Sharing Pathways Directly":** With microservices, the warehouse transforms into a dynamic, collaborative "brain" at the edge. Robot A, specialized in visual recognition, exposes its “High-Accuracy Object Recognition” microservice like offering a highly specialized neural pathway to its neighboring "neurons" (robots). Robot B, the navigation expert, dynamically discovers and utilizes this visual recognition pathway from Robot A, instantly augmenting its own perception capabilities directly, at the device level. In turn, Robot B might share its “Super-Efficient Path Planning” microservice, which Robot A can now leverage to improve its own movement. Robot C contributes its “Inter-Robot Communication” microservice, allowing seamless coordination across the swarm, with all interactions happening device-to-device, at the extreme edge.
- 3) **Direct Capability Transfer: Decentralized "Knowledge Exchange" at the Context Source:** This capability sharing and knowledge exchange happen directly between the robots at the end device edge, mimicking the neuron-to-neuron communication within a brain. No central “warehouse control tower” acting as a bottleneck, no cloud intermediary dictating every interaction. The robots are, in effect, "teaching" and "learning" from each other in real-time,

building a more robust and intelligent collective intelligence within the warehouse, grounded in the rich context of their immediate environment.

This microservice-driven collaboration, happening directly at the end device, unlocks a world of possibilities, allowing for:

- 1) **Dynamic Capability Augmentation: "Borrowing Brain Functions Locally"**: A device, like a single neuron seeking to enhance its functionality, can continuously scan its local network "neighborhood" for available microservices for specialized "brain functions" it might need. If a resource-constrained device encounters a complex task, it can dynamically discover and utilize microservices from nearby devices or edge servers, effectively "borrowing" processing power or specialized AI algorithms right at the edge. When the need subsides or the connection fluctuates, it seamlessly reverts to its own core "neural circuits", maintaining robust local operation.
- 2) **Evolving Edge Ecosystems: "The Brain That Grows Organically at the Edge"**: Microservice architectures foster the creation of dynamic and evolving edge ecosystems, growing organically from the end devices upwards. Imagine a "brain" that is constantly learning and adapting, not through centralized control, but through decentralized interactions and knowledge exchange at the device level. New microservices, representing new skills or knowledge modules, can be dynamically deployed, updated, and shared across devices at the edge, constantly evolving the collective intelligence of the edge ecosystem, and mirroring the plasticity and lifelong learning capacity of the human brain, but now distributed across a network of intelligent devices at the edge of context.

Within the Device-First continuum, end devices are not isolated islands of intelligence but active participants in a collaborative network. Edge-native microservices enable devices to dynamically discover, communicate with, and assist each other, as well as leverage the capabilities of gateways and cloud resources in a coordinated manner. This creates a "social AI dynamic continuum" where devices, gateways, and cloud components work in concert, sharing data, models, and insights to achieve collective intelligence that surpasses the limitations of siloed, cloud-centric, or gateway-bound approaches.

Consider the power of a human team versus a single individual. A well-coordinated team can tackle projects of immense scale and complexity, leveraging the diverse skills and perspectives of its members. Similarly, a DFC-AI system can address challenges that are beyond the reach of isolated AI agents by harnessing the collective intelligence of a vast network of interconnected devices, working in harmony across the flexible continuum. This collaborative intelligence unlocks new possibilities for innovation, efficiency, and human-like problem-solving in a distributed AI world.

## 10 End-to-End Latency Reduction: Achieving Human-Like Reflexes at the Device Edge

Latency is even more critical when we consider the end device as the edge of action. It's not just about reducing network delays. It's about achieving near-instantaneous responsiveness at the point of interaction with the physical world. The Device-First architecture, by running AI agents and much of the relevant workflows on the end device, is uniquely positioned to minimize end-to-end latency and achieve human-like reflexes at the edge of context.

Consider the decision-making pathway in different architectures, imagining a robot needing to react to an obstacle:

- 1) **Cloud-Centric: The "Cloud Reflex": Slow and Distant:** Sensor data from the robot travels across the network all the way to the cloud -> AI processing happens in the cloud -> Decision or Action signal travels all the way back to the robot -> Robot finally reacts. The result is high latency due to extensive network round-trips and centralized cloud processing, like a delayed, conscious decision rather than a reflex.
- 2) **Gateway-Mediated: The "Gateway Reaction":** Faster, but Still a Hop Away: Sensor data travels to a gateway -> AI processing at the gateway -> Decision/Action signal travels back to the robot (or sometimes even a longer path via cloud, then back to the robot). Reduced latency compared to cloud-centric, but still intermediary network hops via the gateway, delaying true real-time response.
- 3) **Device-First: The "Device Reflex" Instantaneous and Local:** Sensor data processed directly on the robot itself -> AI processing happens onboard the device -> Action taken immediately. Minimal latency as the decision is made directly at the source of context and action, achieving near-instantaneous reflexes, truly emulating human-like reaction times.

This latency difference is not just a technical detail. It's a fundamental differentiator for applications where split-second timing is paramount, such as real-time control systems, autonomous vehicles, and many safety-critical systems. It's important to note that for the most demanding safety-critical functions requiring absolute minimal latency and deterministic responses, integration at lower levels, potentially even below the operating system kernel, may be necessary. However, for a vast range of safety-relevant applications and particularly for the overall intelligence, coordination, and higher-level decision-making within even safety-critical systems, the latency reductions offered by DFC-AI provide immense and crucial advantages. In these scenarios, even milliseconds of latency can be the difference between smooth, intelligent operation and system failure or between a human-like reflex and a sluggish, inadequate response.

## 11 Significant Cloud Hosting Cost Reduction: Optimizing Resource Usage like the Human Brain

Let's shift our focus to economics and cost, another area, where DFC-AI, especially when leveraging microservices, mirrors the efficiency of human intelligence. The Pareto principle (80/20 rule) is remarkably applicable to edge AI workloads and their cost implications. Just like our brains efficiently handle the vast majority of everyday tasks locally, minimizing energy expenditure, mimik's platform excels at optimizing cloud hosting costs.

In many hyper-local, context-rich applications, we observe a Pareto distribution:

- 1) **The 80% Local, Context-Dependent Tasks: "Everyday Thinking"**: More than 80% of AI tasks at the edge are relatively small, localized, and highly dependent on immediate context [18], [19]. Examples: simple object recognition (is that a package or an obstacle?), anomaly alerts (is this sensor reading normal?), basic environmental adjustments (is the temperature too high?), local control loops (maintain speed and trajectory). These are like our everyday thoughts and actions handled efficiently and locally by our brains without constant "third-party consultation".
- 2) **The 20% Complex, Global Tasks: "Global Analysis and Learning"**: Less than 20% of tasks are truly complex, require vast aggregated datasets, or benefit from centralized, cloud-based processing. Examples: large-scale model retraining on data from millions of devices, complex analytics across the entire fleet to identify global patterns, infrequent but computationally intensive tasks like major software updates. These are analogous to our more deliberate, analytical, and knowledge-building "global thinking" processes, which benefit from access to broader information and resources. DFC-AI is designed to handle the vast majority (80%+) of these smaller, context-driven, "everyday thinking" tasks directly on the end device, without constant cloud involvement, and delegating the rest to the rest of the cloud continuum. This dramatically reduces the need for continuous, high-volume data transfer and cloud processing, resulting in significant cost savings.
- 3) **Quantifiable Cost Savings: Human-Scale Efficiency, Machine-Scale Savings**: By processing 80%-95% of workloads locally, the DFC-AI approach can demonstrably cut down cloud hosting costs by an average of 80% and up to 95% for many common hyper-local applications [20]. This translates to immense savings in infrastructure expenses, bandwidth consumption, and the vast energy footprint of data centers truly mimicking the economic efficiency of the human brain in handling everyday tasks locally and efficiently.

## 12 Energy Efficiency and Economic Synergies: A Virtuous Cycle of Efficiency Inspired by Nature

The energy efficiency gains of microservices (emulating the brain's specialized circuits) powerfully compound the economic benefits of DFC-AI. By processing

workloads locally and doing so in an energy-optimized way through microservices, we achieve a virtuous cycle of efficiency: reduced energy consumption at both the end device level (longer battery life, lower device operating costs like a brain that doesn't overheat) and significantly reduced cloud hosting bills (less demand on energy-hungry data centers).

This dual efficiency is not just environmentally responsible, but also fundamentally economically sound. It's a path towards building intelligent systems that are not only powerful and responsive, but also sustainable, scalable, and practically deployable in the real world; systems that truly mimic the elegant efficiency and effectiveness of human intelligence, right at the edge.

### **13 DFC-AI's Offline-First, Collaborative, Power-Efficient, and Economically Optimized Agents Embodying Human Intelligence in Machines**

The most compelling vision for DFC-AI is now fully illuminated. It's one that embraces Device-First principles, real-time offline operation, collaborative microservice agents, energy efficiency, and economic optimization, a holistic approach that emulates the very best aspects of human intelligence in our machine creations. This is not just about incremental improvement but a fundamental paradigm shift in how we design and deploy intelligent systems.

### **14 The Need for Context-Aware Agentic Workloads and Collaborative, Efficient, Human-Inspired Intelligence: Demanding a New Era of Edge Intelligence**

The increasing demand for context-aware agentic workloads isn't just a technological trend. It's a reflection of our growing desire for machines that can interact with the world in more sophisticated, adaptive, and yes, human-like ways. Think again of the complex behaviors of humans in dynamic environments: our ability to navigate crowded streets, collaborate in teams, react instantly to unexpected events, all while conserving energy and making economically sound decisions. This is the minimum level of intelligence we are now asking of our end devices.

Just as human intelligence is inherently:

- 1) Context-Aware: Deeply attuned to immediate surroundings and relevant information.
- 2) Agentic: Capable of autonomous decision-making and action.
- 3) Collaborative: Naturally inclined to share knowledge and work together.
- 4) Efficient: Optimized for energy usage and resourcefulness.
- 5) Economically Sound: Operating within real-world constraints and cost considerations.

So too must be our Hybrid Edge AI systems with Device-First architecture fully leveraging all the elements within the AI Continuum that offer the most promising blueprint for achieving this human-inspired level of system-level intelligence.

## 15 Conclusion: Embracing the Flexible AI Continuum for a Human-Like AI Future

In conclusion, the journey towards truly human-like artificial intelligence requires a shift beyond siloed, cloud-centric AI models towards a more distributed, efficient, and collaborative paradigm. DFC-AI offers a compelling path forward. By prioritizing intelligence at the end devices, enabling seamless continuum integration, and fostering collaborative intelligence, it unlocks significant advantages in operational costs, energy efficiency, latency reduction, contextual awareness, and operational effectiveness, especially for real-time, context-aware applications.

As we move towards a future where AI is deeply embedded in our daily lives, from smart environments to autonomous systems, the Flexible AI Continuum enabled by DFC-AI becomes not just a technological advantage, but a necessity. Embracing this paradigm shift is crucial for realizing the full potential of Hybrid Edge AI to create truly intelligent, responsive, and sustainable AI solutions that enhance human capabilities and shape a more intelligent and interconnected world. DFC-AI is not just an architecture; it is a foundation for a human-like AI future.

## References

1. Rob van der Meulen, "What Edge Computing Means for Infrastructure and Operations Leaders", Gartner Technology Insights, October 03, 2018
2. Glikson, A., Nastic, S., Dustdar, S.: Deviceless edge computing: extending serverless computing to the edge of the network. In: 10th ACM International Systems and Storage Conference (SYSTOR), Haifa, Israel, May 2017
3. Shi, W., Cao, J., Zhang, Q., Li, Y., & Xu, L. (2016). Edge computing: Vision and challenges. *IEEE Internet of Things Journal*, 3(5), 637-646. (Broad overview of Edge Computing vision)
4. "What is Cloud Computing?". Amazon Web Services. 2013-03-19. Retrieved 2013-03-20
5. Siavash M. Alamouti, Fay Arjomandi, and Michel Burger. Hybrid edge cloud: A pragmatic approach for decentralized cloud computing. *IEEE Communications Magazine*, 60(9):16–29, 2022
6. Siavash Alamouti, Fay Arjomandi, Michel Burger, Bashar Altakrouri. Building Blocks to Empower Cognitive Internet with Hybrid Edge Cloud. arXiv:2402.00876 [cs.NI, cs.AI]
7. Koch, C. (2004). *Biophysics of computation: Information processing in single neurons*. Oxford University Press.
8. Bassett, D. S., & Gazzaniga, M. S. (2011). Understanding human brain networks. *Nature Neuroscience*, 14(6), 647-654.
9. Barrett, L. F. (2017). *How emotions are made: The secret life of the brain*. Houghton Mifflin Harcourt.
10. Pascual-Leone, A., Amedi, A., Fregni, F., & Merabet, L. B. (2005). The plastic human brain cortex. *Annual Review of Neuroscience*, 28, 377-401.

14 S. M. Alamouti, et al.

11. Bonomi, F., Milito, R., Natarajan, P., & Zhu, J. (2012). Fog computing and its applications. In Workshops at the International Conference on Cyber-Physical Systems (ICCPS).
12. Yi, S., Li, C., & Li, Q. (2015). Fog computing: platform and applications. In 2015 Third IEEE International Conference on Mobile Cloud Computing, Services, and Engineering (pp. 84-89)
13. Laughlin, S. B., & Sejnowski, T. J. (2003). Communication in neuronal networks. *Science*, 301(5640), 1870-1874.
14. Lian Ping Chen, "Microservices: Architecting for Continuous Delivery and DevOps", The Proceeding of the IEEE International Conference on Software Architecture, ICISA 201
15. Baldini, I., Castro, P., Chang, K., Cheng, P., Fink, S., Ishakian, V., Mitchell, N., Muthusamy, V., Rabbah, R., Slominski, A., Suter, P.: Serverless Computing: Current Trends and Open Problems. arXiv:1706.03178, June 2017
16. The serverless trilemma: function composition for serverless computing Baldini, P Cheng, SJ Fink, N Mitchel Proceedings of the 2017 - dl.acm.org
17. Woolley, A. W., Chabris, C. F., Pentland, A., Hashmi, N., & Malone, T. W. (2010). Evidence for a collective intelligence factor in the performance of human groups. *Science*, 330(6004), 686-688.
18. Jackson, K., & Bunch, C. (2019). Understanding Cloud Workloads: An Overview of AWS Instance Types and Usage Patterns.
19. Xu, Y., Bi, J., Yuan, W., & Yuan, W. (2017). Cloud Workload Characteristics and Resource Management Approaches: A Comprehensive Study. *IEEE Transactions on Cloud Computing*, 5(4), 772–785.
20. Siavash Alamouti. Quantifying Energy and Cost Benefits of Hybrid Edge Cloud: Analysis of Traditional and Agentic Workloads. arXiv:2501.14823 [cs.DC, cs.AI]