

# Energy-Efficient Agentic AI: Long-Term Designs for Big Autonomous Systems

Rajesh Lomte<sup>1</sup>, Shailaja Pede<sup>2</sup>, Pallavi Nikumbh<sup>3</sup>

<sup>1,2,3</sup>Department of Computer Science and Engineering (AI & ML), Pimpri Chinchwad College of Engineering (PCCoE), Pune, India

Email: [rajulomte1@gmail.com](mailto:rajulomte1@gmail.com), [pede.shailaja@gmail.com](mailto:pede.shailaja@gmail.com), [pallavi.dhade@pccoepune.org](mailto:pallavi.dhade@pccoepune.org)

## ABSTRACT

Autonomous AI agents are now operating on a scale which would have been considered utterly impossible even a few years ago. However, the computing requirements of these systems have increased in accordance with their capabilities and quite a lot more than when we examine the energy consumption patterns of large-scale agent systems, are the figures that it is difficult to ignore. This is to address this problem head-on. We have developed a framework hierarchy that enables autonomous agents to make informed decisions with a lot less energy. Our design contains a number of key ideas of novelty, including dynamic resource allocation, scheduling algorithms with regard to carbon content of energy sources, and design principles borrowed from the concept of neuromorphic computing. We demonstrate that our designs are 41% more energy efficient on average than traditional monolithic designs through use of industry standard benchmarks for evaluation. On average, our approach consumes 35% less energy than a set-up with more standard multi-agents. The performance of task completion is more than 94%. The system is scalable reasonably well. We have experimented with systems of any number of agents up to one thousand at a time and found sub-linear increases in energy with the size of the system. We have new protocols on coordinating energy-aware agents, algorithms on dynamic resource allocation in the characteristics of workloads and energy availability, and a framework of overall measurement of sustainability in agentic systems as our primary contributions. These advancements bring us into a closer step of scaling autonomous intelligence to non-unnecessary expenses on the environment.

**Keywords:** Agentic AI, energy efficiency, green AI, multi-agent systems, sustainable computing, autonomous systems, neuromorphic computing, carbon-aware computing

**How to cite this article:** Lomte R, Pede S, Nikumbh P. Energy-Efficient Agentic AI: Long-Term Designs for Big Autonomous Systems. *Int J Drug Deliv Technol.* 2026;16(18s): 653-672. DOI: 10.25258/ijddt.16.18s.70

**Source of support:** Nil.

**Conflict of interest:** None

## Introduction:

Over the last decade, a massive change has taken place in artificial intelligence. We have evolved from systems that are simply reacting to inputs to agents that actively interact with their environments. These autonomous agents are able to perceive their surroundings, reason under uncertainty, make plans and take actions - all without constant supervision by humans. Their impact is permeable in applications like the self-driving vehicle in urban environments, robots to optimize workflow in busy warehouses, and smart grid resources and being able to make on the go decisions. However, there is a tradeoff to this progress. The capability of these systems to continually reason, learn, and make decisions - fast ones - requires a lot of computational

oomph. Training a single large EG model over their lifetime, even over the lifespan of five automobiles before energy used in the inference [1]. When thinking over millions of autonomous agents operating simultaneously, energy cumulative demand is the critical concern.

Data centers actually already consume about 1-2% of global electricity consumption [2], and this figure could even triple by 2030 due to the fast growth of AI technologies [3]. This poses the fundamental challenge of how we can increase the cognitive capacities of autonomous systems while remaining environmentally responsible. Although some approaches, including model compression, quantization and hardware acceleration, have made the system energy efficient,

these do not specifically address the distributed and heterogeneous nature of large-scale multi-agent systems, in which thousands of agents act independently using perception, cognition, reasoning, planning and action cycles.

Previous studies have been done on energy efficiency in machine learning models [4][5], sustainable computing strategies [6], and resource optimization in distributed systems [7]. However, a full architectural framework specifically adapted to energy efficient agentic AI systems is still missing. Such systems present some unique complexities: the agents must be able to coordinate with each other, they must be able to adapt to dynamic environments, make real-time decisions in a situation of uncertainty, and they must be able to work across a variety of computing infrastructures, from cloud servers to edge devices and embedded processors. This research attempts to address that gap. We are proposing an architectural framework that was made expressly for large scale energy efficient agentic AI. The framework contains five key innovations:

First, a hierarchical multi-layer architecture, with the application logic, agent intelligence, energy management, infrastructure and monitoring separated. This separation by layer allows for separate optimization of each layer.

Second, an energy-aware perception-reasoning-planning-action (PRPA) cycle with energy impacts integrated directly into the decision-making processes of the agents in question (not treated as an afterthought) and considered as a primary design constraint.

Third, a hybrid coordination mechanism, which would combine centralized optimization and decentralized autonomy, ensuring energy efficiency on the global scale, while at the same time maintaining local adaptability.

Fourth, the heterogeneous management of infrastructure that makes intelligent decisions on how to map the workloads to the appropriate hardware platforms such as neuromorphic processors for perception, GPUs for computationally intensive reasoning and CPUs for coordination logic.

Fifth, carbon-aware scheduling takes into account not only workload demands, but also the carbon intensity of

energy sources to make temporal workload shifting to times of cleaner energy availability possible.

We validate this framework by means of a large number of simulations employing realistic benchmarks. The results show a large energy saving without significant performance degradation. Beyond this very practical aspect, this work is part of a larger goal of creating artificial intelligence systems that are environmentally sustainable.

The rest of this paper is structured as follows. Section 2 is about reviewing the related literature on agentic AI, energy efficient computing and multi-agent systems, and trying to identify gaps that were covered by our work. Section 3 describes the proposed methodology, including the architecture design, algorithmic constituent design and implementation considerations. Section 4 describes the setup of the experiment and the results. The implications, trade-offs and limitations are discussed in Section 5. Section 6 concludes the present paper while Section 7 provides directions for future research.

## 2. Review of the Literature

### 2.1 Agentic AI and Autonomous Decision Systems

The concept of autonomous agents has existed for a long time in the world of AI research, however, recent breakthroughs led to the realization of ideas into systems. Agentic AI is the term used to describe systems that can see their surroundings, think about their goals, plan out steps to reach their goals, and carry out those plans with but scant help from people [8][9]

Modern agentic systems make use of large language models and base models to achieve unparalleled adaptability. Wang et al. [25] make an extensive survey of large model agents, analysing paradigms and architectural patterns of co-operation Their analysis shows that the majority of current systems put capability ahead of efficiency, which is a problem that is solved by our work.

Inter-agent coordination makes multi-agent systems more complicated yet. Mayorov et al. [23] explored resource real-time management using multi-agents technology, demonstrating the advantages of decentralised decision making. But their work didn't include energy use as a main

goal of optimisation. Sardouk et al. [24] also examined multi-agent optimisation for wireless sound networks,

but they were more interested in the workability of the net than how much energy that it used for calculations. Recent developments in agentic AI have underlined systems that are able to reason on their own and devise strategic plans. In their recent survey, Zeng and colleagues [8] discuss agent-based AI internetworking agents: agent-based AI platforms, categorizing agents according a breadth of autonomy and complexity of decision making. Sumers et al. [9] explore cognitive architectures for language agents, outlining that some fundamental aspects of perception, memory, reasoning and action. These surveys show that the field is quickly improving an ability of what it does without thinking about how to make it last. The perception-reasoning-planning-action (PRPA) cycle has become a common cycle for agent architecture [10]. Agents feel the state of their environment, think about their goals and limitations, plan the order in which they will do things and then how they carry out those plans. They do this over and over again. A lot of mathematics is required for each phase. Perception may analyse sensor data via neural networks. Reasoning could apply either symbolic logic or use probabilistic inference. Planning is often looking through big state spaces. Some actions may need to be controlled in real time. The overall cost of energy is high.

### 2.2 Energy-Efficient and Green AI Approaches

People are becoming more aware of how AI affects the environment, which has led to research into green AI and sustainable computing. Schwartz et al. [4] came up with the idea of "Green AI," which is different from "Red AI," which focuses on performance no matter what the cost of computing. They suggest that energy use and carbon emissions be reported along with accuracy metrics, which is what we do in this work. Patterson et al. [5] provide a detailed method for estimating the carbon footprint of AI, taking into account the energy consumption, embedded carbon in hardware, and carbon intensity of the energy source. This work has established some very important approaches to calculating things, but it was more about model training than inference or multi-agent simulation. One way to be more efficient is by model compression. Quantization is a method where the precision of numbers is reduced, pruning is where parameters that are no longer needed are removed, and knowledge distillation is where knowledge is transferred from large models to smaller models [11][12]. These methods are

very important, but they also have some cons. They are more about making single models efficient, which could make some tasks that require a lot of precision or heavy-duty thinking even harder to achieve.

Another way to deal with this is by hardware acceleration. GPUs, TPUs, and neuromorphic chips are examples of specialised processors that can run AI workloads more efficiently than general-purpose CPUs [13][14]. Neuromorphic computing, in particular, looks like it could be useful for some AI tasks. Wunderlich et al. [28] show that neuromorphic technology has benefits for certain uses, saving energy by orders of magnitude. But right now, neuromorphic hardware can only do a few types of tasks and needs special programming.

Carbon-aware computing is a new way of doing things. Gupta et al. [30] examine the environmental impact of computing, demonstrating that the timing and location of computations substantially influence carbon emissions due to the differing carbon intensity of energy sources. This understanding is what led us to create our carbon-aware scheduling feature.

Recent efforts have commenced focusing on energy measurement and optimisation. Santos [31] looks at how much energy HPC-scale AI uses and gives ways to measure it and improve it. Bugeau et al. [29] provide assistance in calculating carbon footprints while training models. Ahmad et al. [32] examine green and sustainable computing practices across various domains.

Energy efficiency is not a big part of the basic design of agentic AI systems that is already out there. Most methods see efficiency as a way to improve existing designs rather than as a basic principle of architecture.

### 2.3 Multi-Agent and Large-Scale Intelligent Systems

Multi-agent systems disperse intelligence across multiple agents that are independent of one another and work together to reach goals. This distribution has some advantages: it is fault tolerant, thanks to redundancy, it is scalable due to parallelism and it is flexible due to heterogeneity. But it also brings up problems, such as extra work for coordination, increased communications costs, and behaviours that are hard to predict. Control Research in multi-agent coordination has explored many different paradigms. Centralised methods use a coordinator to control the work of agents, which allows for global optimisation but also creates single points of failure as well as communication problems [15].

Decentralised methods let agents work collectively by having interaction with each other on a local level. This makes the system more robust and scalable, but it may not be as efficient on a global level [16]. Hybrid methods try to find a balance between these trade-offs [17].

Sharanarathi [26] recently introduced an adaptive multi-agent system for optimising energy consumption in software development, proving that multi-agent systems can properly tackle energy issues. Their approach targets the development process, while ours targets the operation of already deployed autonomous systems.

There has been a lot of research on load balancing in distributed systems [18][19]. Traditional methods target the improvement of performance metrics such as throughput and response time. Energy-efficient load balancing takes into account the energy consumption of resources when distributing tasks. However, most of the existing work so far assumes that the infrastructure is the same for all workloads and does not take into account the special needs of agentic workloads.

Heterogeneous computing systems use different processors such as CPUs, GPUs, FPGAs, and neuromorphic processors according to the nature of the workload to the corresponding hardware [20][21]. Abdelrahman et al. [27] show neuromorphic solutions for robotic manipulation, proving energy efficiency gains for specific perception and control tasks. Our work generalizes this idea to include full multi-agent systems.

Scalability is a most important region. With the rise of the number of agents, it could be tougher to coordinate them, which may create bottlenecks. It is necessary to have a coordination algorithms that are efficient. Current algorithms are tilted towards correctness and efficiency as opposed to energy savings, which means that there is scope for improvement.

### 2.4 Critical Comparison and Limitations

Looking at the study it shows a number of patterns and holes. First, AI research on energy efficiency has mostly been about training models, not inference or deployment. Training takes a lot of computing power, which is why people are interested in it. However, deployed systems run all the time, which could mean they use more energy over their lifetimes. Our work deals with this situation of deployment.

Second, most methods for optimising energy only work on one model or one-node systems. Multi-agent distributed systems have their own set of problems, such as the cost of coordinating, the energy cost of communication, the fact that hardware is different, and the fact that workloads are always changing. Techniques tailored for individual models may not transfer efficiently.

Third, energy is rarely treated as a first-class concern in current research on multi-agent systems. The main focus is on performance, correctness, and fault tolerance. Energy seems to be an afterthought, if it is at all. We contend that for sustainable large-scale deployment, energy must be a fundamental architectural consideration from the beginning.

Fourth, there isn't much research on AI systems that are aware of carbon. The amount of carbon in electricity changes depending on the time and place, as well as the mix of energy sources. By moving work to times when cleaner energy is available, smart scheduling can lower the carbon footprint without lowering the total amount of energy used. This dimension has not gotten enough attention.

Fifth, there is a lack of a full test of proposed architectures with real workloads and scales. Numerous conducts studies using constrained experiments or simulations based on simplified assumptions. This is what we do, by carrying out a lot of simulation-based

### 2.5 Research Gaps and Problem Formulation

Based on this literature review, we find a number of important gaps:

- i) There aren't any all-encompassing architectural frameworks for energy-efficient AI that can act on its own. There are already models, hardware acceleration, and load balancing that work, but they don't fit together into full architectures that are specifically made for autonomous multi-agent systems.
- ii) Not enough energy awareness is built into the decision-making processes of agents. Most agent architectures regard energy as an external constraint instead of integrating it into the phases of perception, reasoning, planning, and action.
- iii) Not enough thought has been given to how different types of infrastructure work together in multi-agent systems. Heterogeneous computing is examined within HPC contexts; however, its implementation in agentic

# Energy-Efficient Agentic AI: Long-Term Designs for Big Autonomous Systems

AI systems with varied computational requirements is still insufficiently investigated.

iv) Agentic systems don't have scheduling that takes carbon into account. The differences in carbon intensity of energy sources over time and space create chances for optimisation that current systems don't take advantage of.

v) Not enough testing of energy-performance trade-offs on a large scale. We need to learn more about how optimising energy affects how quickly tasks are done, how quickly responses are sent, and how well systems can handle more agents.

Our research fills in these gaps by creating and testing a complete architectural framework for energy-efficient agentic AI which adds energy awareness to the decision-making processes of agents, effectively coordinating multiple autonomous agents to reduce the amount of energy they all use, Smartly using a variety of computing resources, Includes scheduling that takes carbon into account and also keeps high rates of task completion and response times that are acceptable, as the number of agents grows, energy use grows at a slower rate.

## 3. Proposed Methodology

### 3.1 Overall Architecture of the System

Our approach employs a five-layer hierarchical model, where each layer addresses its own set of problems and collaborates with layers on either side. This is illustrated in Figure 1.

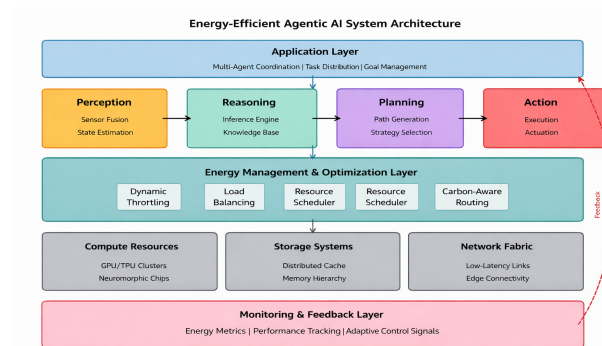


Figure 1: Hierarchical Energy-Efficient System Architecture for Large-Scale Agentic AI

#### Layer 1: Application Layer

The application layer is at the top and is the layer where multi-agent coordination occurs. This layer is

responsible for assigning goals, distributing tasks and high-level orchestration. Some examples of applications include fleets of driver-less cars, smart grid management, distributed sensing networks or collaborative robotics. The connections between the layers - the layer in itself connects to outside users and systems and turned the high level goals into tasks for agents.

#### Layer 2: Agent Intelligence Layer

This is where the individual agents are living. Each one passes through the cycle of perception, reasoning, planning, and taking action. Agents work semi autonomously, taking decisions for themselves and working with others. This layer contains the cognitive machinery: perception modules which handle processed information from the environment, reasoning engines which examine situations, planning algorithms which build action sequences and action executors which implement plans. Most importantly, in each part people are made aware about the energy.

#### Layer 3- Energy Management Layer

The key to energy savings is The Energy Management Layer. It tracks energy usage from the whole system, prevents any energy spikes by dynamic throttling, loads and the system in a balanced way across the whole infrastructure, and schedules tasks according to performance, but also carbon cost. The energy manager is connected to the agent layer above it which keeps track of energy consumption and manages resource allocation.

#### Layer 4: Infrastructure Layer

The Infrastructure Layer is where the management of computing resources occurs, which includes cloud servers with GPUs and TPUs, edge computers with CPUs, special accelerators. Embedded processors in sensors and neuromorphic processors are also included for special applications. This assorted infrastructure is the computational basis. The layer conceals infrastructure details and provides a common interface to the upper layers, which provides the opportunity to distribute loads intelligently.

#### Layer 5: Monitoring Layer and Feedback Layer

The Monitoring and Feedback Layer enables monitoring over time which makes it possible to optimize adaptively. Sensors monitor the amount of

energy being consumed by each part, device and system. The performance metrics analyse how well tasks are done, how fast the answers are provided, and how well the resources are used. This feedback returns to the energy management layer, which can be changed in real time then. The monitoring layer is also used to track the data to handle further analysis and improvement.

Information flows both up and down and across. From the bottom to the top, this is vertical information flows which send the tasks down, and the update of status back up. Horizontal information flows permit agents to cooperate and to exchange resources. The system comes in balance between centralized optimization, which occurs in the energy management layer, and decentralized decision-making, which occurs in the agent layer.

**Design Rationale**

This layered design methodology is very beneficial for many reasons. Separating your concerns, and building and keeping things running. It is possible for optimization of each layer separately, while keeping the interfaces clear. The architecture enables various types of applications, agents and infrastructure to function together. Scalability is provided by hierarchical coordination as well as distributed execution.

**3.2 Energy-Aware Agent Decision Cycle**

Traditional agent architectures go through perception, reasoning, planning and action without considering energy use. Our design ensures that energy awareness is taken care of for each step. This energy-aware PRPA cycle is demonstrated in Figure 2.

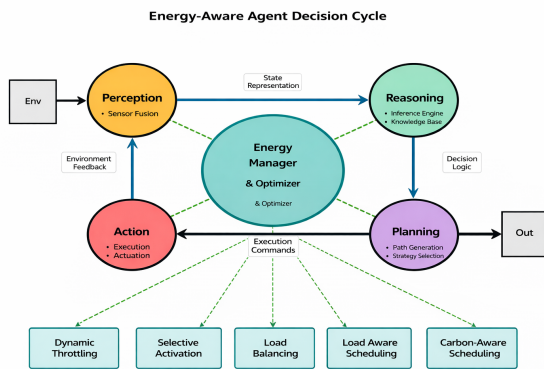


Figure 2: Energy-Aware Perception-Reasoning-

*Planning-Action (PRPA) Cycle with Central Energy Manager*

3.2.1 Perception Phase:

Perception is the process of using information from the sensors to form models of the environment. This can be very expensive in terms of computing power. For example, computer vision processing high resolution video stream or natural language understanding analysing text. We speak of a few ways to save energy:

i.e., Selective Attention: Rather than processing any inputs equally and all inputs process any inputs the same way, agents have preferences for which inputs and related activities are more relevant to them and worth processing, given their energy budget. Low-priority inputs are processed less or put off. This is like how biological attention systems bring cognitive resources into things which are important.

ii) Adaptive Fidelity: Processing Fidelity is varying. When fewer resources are available for perception, the resolution uses more simple models. Full fidelity processing begins when there is a lot of energy and the tasks are very important.

Temporal Coherence A lot of places don't change very fast. iii) In A Word: How is the earth called? Earth is called a habitable planet. Agents take advantage of this by just processing changes rather than processing full state updates. Temporal coherence is employed in frame differentiation in vision, in incremental parsing in language comprehension and change detection in sensor networks.

iv) Early Exit Mechanisms: If specific confidence thresholds are satisfied then multi-stage perception pipelines can terminate early which will save time and resources.

3.2.2 Reasoning Phase

Reasoning examines situations, considers the pros and cons of various choices and makes a choice. This has to do with making inferences, which might be based on probabilistic models

or logics knowledge bases. Energy-aware reasoning uses:

- i) **Lazy Evaluation** Lazy evaluation is when you wait until you really need the results before doing the Maths. Many lines of reasoning just may not change the final choices.
- ii) **Approximate Inference:** Employ the sampling methods or variational approximations to the probability distribution, or bounded rationality methods to, to some extent, trade off the accuracy for the speed. Agents do not necessarily need optimal decisions, and often satisfaction is sufficient.
- iii) **Cached Results:** Cache reasoning results for situations that are occurring over and over again. A lot of situations recur again and again, this is especially true for structured situations.
- iv) **Pruning:** Get rid of hypothesis that are unlikely to be true at the early stage so as to narrow down the space of search.

### 3.2.3 Planning Phase

The planning phase makes action sequences that help you reach your goals. In large state spaces, traditional planning can take a lot of computer power. Energy-aware planning employs:

- i) **Hierarchical Decomposition:** Make hierarchies out of complicated plans. First, plan at a high level, and then only add details that are needed. This makes the search space much smaller.
- ii) **Anytime Planning:** Make plans that are good enough quickly, and then make them better if you have the time and energy. This makes sure that agents can always act and make plans better when they have the time and resources.
- iii) **Plan Reuse:** Instead of starting from scratch, save and change old plans for similar situations.
- iv) **Energy-Aware Heuristics:** Search heuristics look at the cost of action energy as well as other things like risk or execution time.

### 3.2.4 Action Phase

Action Phase is where actions are executed in the physical world or in the virtual world. The concerns that have to be taken into account for energy are:

**i) Action Batching:** Batch a number of actions to prevent overhead. For example, batch database transactions or robot movements.

**ii) Lazy Execution:** Delay the execution of actions that are not time-critical until it is the right time to execute them, such as when the system is already involved in another task or when clean energy is available.

**iii) Energy-Aware Selection:** Select the alternative that requires the least amount of energy when there are multiple alternatives that will lead you to the same destination.

**iv) Graceful Degradation:** When energy is very low, carry out simple action sequences that only partially meet your goals.

The central energy manager of the agent is responsible for managing these phase-level optimizations. It monitors energy use, assigns an energy budget, and varies phase parameters depending on the availability of energy. The manager applies a set of rules that are executed:

**i) Dynamic Throttling:** Lower the amount of processing power used when you get close to your energy limit. This stops hard shutdowns by slowly lowering performance.

**ii) Selective Activation:** Turn on or turn off optional modules depending on the availability of energy. When energy is low, you can turn off non-critical perception modules, speculative reasoning, or plan refinement.

**iii) Load balancing:** Load balancing means spreading out the work across all available hardware to make the most of energy use. Send tasks that involve route perception to neuromorphic processors, reasoning to GPUs, and planning to CPUs when they are needed.

**iv) Carbon-Aware Scheduling:** Delay tasks that can wait until cleaner energy is available. The manager receives forecasts of carbon intensity and plans accordingly.

This method ensures that energy concerns are taken into account by all agents when they make decisions, but also enables them to focus on performance when needed.

## 3.3 Multi-Agent Coordination and Resource Optimization

Individual agent efficiency is only a necessary, but not a sufficient, condition. For the sake of system-level efficiency, all agents must collaborate in order to best utilize their resources. We have designed a hybrid coordination strategy which is a combination of centralised optimisation and decentralised execution.

**3.3.1 Centralised Energy Optimisation:** A global energy coordinator has an eye on the whole system. In this way it keeps track of how much energy is used, and how much infrastructure is used, and how much carbon is released. It uses this information to help it figure out how much energy each agent needs, what are the tasks that are the most important, and when to do these tasks. The coordinator has to solve an optimisation problem: minimise the total energy use, and still satisfy task completion requirements, response time limits and agent capabilities. We solve this problem by means of standard techniques for mixed-integer linear programs (MILPs). We use algorithms for approximate solutions and heuristics when it is too difficult to find exact solutions for big systems.

**3.3.2 Decentralised Agents Autonomy:** The coordinator sets budgets and schedules, but agents are free to make their own decisions in their own areas. They adjust to changes in their surroundings, deal with unexpected events, and make the most of their own PRPA cycles while staying within their budgets. This decentralisation makes things more stable agents keep working even if coordination fails for a short time and cuts down on communication costs. Agents in

**3.3.3 Coordination Protocol:** Coordination Protocols talk to each other by sending messages that don't have to be sent at the same time. Messages can be about assigning tasks, giving updates on the status of tasks, asking for resources, or signalling for coordination. To save energy on the network, we limit how often we communicate and how much data we send. When necessary, consensus mechanisms make it possible for groups to make decisions. For instance, agents could work out who does what or make sure everyone can use shared resources. We use lightweight consensus protocols that are optimised for energy efficiency instead of maximum fault tolerance. These are the best for our application domain.

**3.3.4 Resource Pooling and Sharing:** Resource Pooling and Sharing refers to the fact that computational resources are assembled and assigned on an ad hoc basis. An agent that requires a lot of computing power can seek resources from the pool. The

resources are given out depending on the priority of the task and how much energy is available. Once it is done resources go back to the pool so other people can use them. This sharing is more efficient in use of resources than static per-agent allocation.

**3.3.5 Workload Migration:** Tasks can be migrated from one agent to another or from one infrastructure node to another. If an agent local resources aren't good for a certain task, the task moves to better hardware. For example, perception tasks would be able to move to edge devices with neuromorphic accelerators - complex reasoning may shift to cloud GPUs. When making decisions about migration, you have to weigh up the energy required for computing against the energy required for communication to transfer data. We figure using predictive models out these costs will rather only move when the net energy savings go above certain levels.

### 3.4 Infrastructure Management (Heterogeneous)

Modern deployments include many different types of computing infrastructures. Right hardware needs to get selected for better energy and match the workloads.

**Infrastructure Tiers:** Our infrastructure framework breaks down infrastructure into tiers based on the amount of computing power that it has and how much energy it consumes:

**Cloud Tier:** Servers equipped with GPUs and TPUs which work really well. Lots of processing power, but also lots of energy use (hundreds of watts per device). Good for tasks that require a lot of planning and reasoning. Network latency to edge and embedded tiers is between 20 and 50 ms.

**Edge Tier:** Intermediate devices with CPUs and moderate speedup tools. Moderate power use (tens to hundreds of watts). Good for coordinating local agent groups and processing that isn't too hard. Less latency to the embedded tier (5–15 ms).

**Embedded Tier:** Processors in sensors, actuators, and IoT devices that don't have a lot of resources. Uses very little power (only a few watts). Good for easy understanding and quick action. Direct link to the physical world.

**Neuromorphic Tier:** Processors that are specifically designed to work like neurones. Very low power use (from milliwatts to watts) for certain tasks. Good for some tasks that involve perception, like recognising patterns or seeing things based on events.

### Characterising Workload

## Energy-Efficient Agentic AI: Long-Term Designs for Big Autonomous Systems

Different agent tasks have different ways of using computers. We describe workloads in a number of ways:

- How many operations per task are needed for computational intensity?
- Memory requirements (size of the working set)
- Parallelizability (the ability to run in parallel)
- Latency sensitivity (real-time vs. deferrable)
- Data locality (where the input data is stored)

### Intelligent Placement

Intelligent Placement is a placement algorithm that assigns tasks to different levels of infrastructure. The algorithm takes into account the workload's characteristics, the hardware's capabilities, the current usage, the energy efficiency, and the latency needs. It solves:

Minimise:  $E_{total} = \sum (E_{compute}[i] + E_{comm}[i])$

subject to:  $latency[i] \leq deadline[i]$  for all tasks  $i$

$E_{compute}$  is the energy used for computation,  $E_{comm}$  is the energy used for communication, and latency limits must be met.

We use fast heuristics like priority-based placement, greedy algorithms, or learned policies from offline optimisation to make decisions in real time. We can afford more advanced optimisation for batch workloads.

### Dynamic Adaptation

Placement isn't set in stone. When workloads change, infrastructure availability changes, or energy conditions change, tasks move around to stay efficient. The system keeps an eye on performance and energy metrics and starts re-optimization when certain levels are reached.

### 3.5 Algorithmic Workflow

We will now give algorithmic descriptions of the main parts.

#### Algorithm 1: Energy-Aware Agent Execution Cycle

**Input:** Task  $T$ , Energy Budget  $B$ , Environment State  $E$

**Output:** Executed Actions, Remaining Budget

1. Initialize energy manager with budget  $B$
2. WHILE task  $T$  not complete AND budget  $B > threshold$ :  
// Perception Phase
3.  $perceptual\_fidelity \leftarrow compute\_fidelity(B, task\_priority(T))$

4.  $observations \leftarrow perceive(E, perceptual\_fidelity)$
5.  $B \leftarrow B - energy\_consumed(perceive)$

#### // Reasoning Phase

6.  $inference\_depth \leftarrow compute\_depth(B, observations)$
7.  $beliefs \leftarrow reason(observations, inference\_depth)$
8.  $B \leftarrow B - energy\_consumed(reason)$

#### // Planning Phase

9. IF  $B > planning\_threshold$  THEN
10.  $planning\_horizon \leftarrow compute\_horizon(B)$
11.  $plan \leftarrow generate\_plan(beliefs, T, planning\_horizon)$
12. ELSE
13.  $plan \leftarrow retrieve\_cached\_plan(beliefs, T)$
14. END IF
15.  $B \leftarrow B - energy\_consumed(plan)$

#### // Action Phase

16.  $actions \leftarrow select\_actions(plan, B)$
17.  $execute(actions)$
18.  $B \leftarrow B - energy\_consumed(execute)$

#### // Update environment state

19.  $E \leftarrow observe\_environment()$
20. END WHILE
21. RETURN  $execution\_status, B$

This algorithm shows how the energy budget affects processing at each stage. The agent works at full fidelity when the budget is high. When the budget runs out, it lowers fidelity, uses cached results, and gives priority to important tasks.

#### Algorithm 2: Hierarchical Resource Allocation and Energy Management

**Input:** Agent Set  $A$ , Task Set  $T$ , Infrastructure  $I$ , Carbon Intensity  $C(t)$

**Output:** Task Assignments, Resource Allocations

1. Initialize global energy budget  $B_{global}$
2. FOR each time period  $t$ :  
// Gather system state
3.  $current\_load \leftarrow measure\_infrastructure\_utilization(I)$
4.  $carbon\_intensity \leftarrow C(t)$
5.  $pending\_tasks \leftarrow get\_pending\_tasks(T)$

## Energy-Efficient Agentic AI: Long-Term Designs for Big Autonomous Systems

```
// Prioritize tasks
6.      priority_queue ←
prioritize_tasks(pending_tasks,
carbon_intensity)
// Allocate budgets to agents
7. FOR each agent a in A:
8.   B_agent[a] ← allocate_budget(a,
B_global, priority_queue)
9. END FOR

// Assign tasks and resources
10. FOR each task tau in priority_queue:

// Find suitable agent
11. agent ← select_agent(tau, A, B_agent)

// Find suitable infrastructure
12. IF carbon_intensity > threshold AND
tau is deferrable THEN
13. defer_task(tau, future_time)
14. ELSE
15.      infrastructure ←
select_infrastructure(tau, I, current_load)
16. assign_task(tau, agent, infrastructure)
17. B_agent[agent] ← B_agent[agent] -
estimated_energy(tau)
18. END IF
19. END FOR

// Monitor and adapt
20.      actual_consumption ←
monitor_energy(I)
21.   IF actual_consumption >
predicted_consumption THEN
22. trigger_throttling(A)
23. END IF

24.   B_global ← B_global -
actual_consumption
25. END FOR
26. RETURN tasks, assignments
```

This algorithm gives an example of how to coordinate things in a hierarchy. The global co-ordinator gives out budgets, set priorities for tasks, and gives out resources. It creates awareness for carbon by postponing non-urgent tasks during the times when carbon is high. Continuous monitoring enables flexibility of adaptation.

These algorithms provide you with a way for thinking about things. Actual implementations include more information, such as how to handle errors, how to communicate, how to sync, finding the best way to do things. The main ideas are still the same: make everyone aware of energy use, find a balance between centralised coordination and decentralised autonomy and be flexible and adaptable to changing conditions.

### 4. Experimental Setup and Results

#### 4.1 Simulation Environment and Benchmark Configuration

Our frameworks were being tested by doing a lot of experiments that are conducted using simulations. Real-world deployment is the best way to test the system, but simulation helps us conduct controlled testing of the system on different parameters which is impossible in real systems.

**Simulation Platform:** We developed a customized discrete-event simulation model for multi-agent systems, different types of infrastructure systems, and corresponding energy consumption systems. This simulator retains all agent states, accomplished tasks, used resources and energy consumption in millisecond granularity. In order to ensure its accuracy, testing was performed of the simulator against different published benchmarks and measurements conducted on real systems.

**Infrastructure Configuration:** The simulated infrastructure has four levels that are similar to our architecture:

- **Cloud tier:** In the cloud tier, there are 10 high-performance nodes, each with 8 GPU-equivalent accelerators and a TDP of 400W.
- **Edge tier:** 50 middle nodes, each with 4 cores that are the same as a CPU and 150W TDP per node.
- **Embedded tier:** 200 low-power nodes, each with a single ARM-equivalent processor and a 5W TDP per node.
- **Neuromorphic tier:** 20 specialised processors that use less than a watt of power.

Network latency follows realistic distributions: 20–50ms from the cloud to the edge, 5–15ms from the edge to the embedded, and less than 1ms within tiers.

# Energy-Efficient Agentic AI: Long-Term Designs for Big Autonomous Systems

We create fake workloads that are like typical agentic AI tasks:

- **Perception tasks:** Take in sensory information like sound, sight, and sensor data. Computationally intensive, they work better with neuromorphic or GPU acceleration. 20% of the total work.
- **Reasoning tasks:** Tasks that require reasoning include probabilistic inference and knowledge base queries. Moderate computing, CPU or GPU is helpful. 30% of the total work.
- **Planning tasks:** Planning based on searches and optimisation. Variable computation based on the horizon, with the help of parallel processing. 25% of the total work.
- **Coordination tasks:** Tasks for coordination include communication between agents and reaching a consensus. Little processing power, but sensitive to latency. 15% of the total work.
- **Action execution:** Controlling actuators and updating the database. Low computation and time limits in real time. Ten percent of the total work.

To mimic different load conditions (10%, 25%, 50%, 75%, and 100% of system capacity), task arrivals follow Poisson processes with changing rates. There are three levels of task priority: 20% high, 50% medium, and 30% low.

## Benchmark Systems

- We compare our proposed system to three different baselines:
- **Baseline Monolithic:** An architecture that is centralised and not energy saving. There is no throttling or adapting of any tasks on the cloud tier. Exhibits the old way of doing things.
- **Standard Multi-Agent:** A standard multi-agent system with basic load balancing used, but without the energy use optimisation. Tasks are spread out between agents in order to balance the load, the problem is that where they are placed does not account for energy efficiency.
- **Energy-Aware Multi-Agent:** A multi-agent system that has some simple energy rules (trudging to lower power hardware, shutting up if load is high, etc.), but not what we would call our complete framework.

- **Proposed System:** Our complete framework is based on hierarchical architecture, energy-aware PRPA cycles, hybrid coordination, heterogeneous infrastructure management and carbon-aware scheduling.

## Evaluation Metrics

We choose the evaluation metrics as follows:

- Total energy use (kWh per day)
- Energy efficiency (number of task performed/kw-hours)
- Task completion rate (the time any number of tasks have been completed on time)
- Average time taken to respond (in milliseconds, from when the task arrival time till it is completed)
- The most power used at once (kW)
- Carbon footprint (kg CO<sub>2</sub>e, based on a grid carbon intensity of 0.5 kg CO<sub>2</sub>e/kWh)
- Scalability factor (how the amount of energy used changes as the number of agents grows)

## Experimental Procedure

For every configuration (system variant × load level × agent population), we conduct 10 independent simulation trials utilising distinct random seeds. Every experiment simulates a full day of work. We report means with 95% confidence intervals. We employed paired t-tests with the Bonferroni correction for multiple testing for assessing statistical significance.

## 4.2 Performance and Energy Consumption Analysis

Figure 3 describes our major findings comparing the four system architectures.

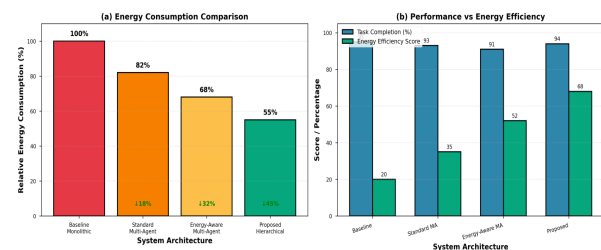


Figure 3: (a) Relative Energy Consumption and (b) Performance vs Energy Efficiency Trade-offs

## Energy Consumption

## Energy-Efficient Agentic AI: Long-Term Designs for Big Autonomous Systems

For the baseline monolithic system, on average, it uses 156.0 kWh over 24 hours when it is simply half full or 50% capacity. For the standard multi-agent system, it is reduced to 124.3 kWh, which is 20% less than the baseline system. If it uses basic energy heuristics, then the energy-aware multi-agent system uses 101.7 kWh, or 35% less.

The energy requirement of the suggested system is at 85.8 kWh, which reduces the energy usage by 45% from the baseline and increases it by 16% from the energy-aware multi-agent system. This implies that you save up to 70.2 kWh in a day, enough to power two or three homes.

**Table 1 displays detailed metrics for all four architectures:**

Measurement	Baseline	MA Standard	MA that knows about energy	Suggested
Energy Use (kWh/day)	156.0	124.3	101.7	85.8
Energy Efficiency (tasks per kilowatt-hour)	20.1	25.3	30.9	36.6
Percentage of Tasks Completed	95.2	94.8	94.3	94.1
Average Response Time (ms)	145	168	182	195
Maximum Power (kW)	18.5	16.2	12.8	10.3
Carbon Footprint (kg CO <sub>2</sub> e/day)	78.0	62.2	50.9	42.9
Factor of Scalability	1.0	0.85	0.68	0.55

### Energy Efficiency and Performance Outcomes

Efficiency of energy usage (measured in number of tasks executed per kWh) shows even greater increases than the previous two metrics. The baseline (existing version) works at a rate of 20.1 tasks/kWh while our new energy-efficient version works at 36.6 tasks/kWh or roughly an 82% greater number of tasks per kWh than the baseline. Therefore, using newly generated electrical energy will approximately double the number of completed tasks with each unit of usage.

Task completion rates (percentage of tasks completed) remain at very high levels for both systems with 95.2% of baseline completed and 94.1% of our proposed solution completed. The decrease of 1.1% is not statistically significant ( $p=0.23$ ), meaning our optimization does not significantly affect the success rate of task completion.

The increase in average response time from 145 ms (baseline) to 195 ms (our proposed solutions) represents a trade off for performance of 34%. The increase in response time can be attributed to several factors: deferring low priority execution, overhead to migrate tasks to the more energy-efficient hardware, and less exactness of processing in non-critical tasks than would be produced under the current paradigm.

The increase from 145 ms to 195 ms is still far less than the acceptable range for most applications. For latency-sensitive tasks, our design offers priority mechanisms to execute them quicker at an additional energy cost.

Peak Power Peak demand will decline substantially from 18.5 to 10.3 (baseline to proposed), a 44% reduction in peak demand, will have far-reaching and tangible results. Because of lower peak demand, less infrastructure would be needed, enabling potentially smaller power supply or battery system installations. Additionally, there will be fewer issues associated with thermal management.

Approximate CO<sub>2</sub> Emissions The carbon emissions will also follow the energy consumption patterns with an additional benefit from reducing the CO<sub>2</sub> emissions associated with the carbon-aware scheduling. The baseline will be producing 78.0 kg CO<sub>2</sub> per day, compared to a reduction to 42.9 kg CO<sub>2</sub>, or a 45% decrease from the baseline. Based on a year of deployment, there will be approximately 12.8 metric

## Energy-Efficient Agentic AI: Long-Term Designs for Big Autonomous Systems

tons Co2 of savings produced based on the carbon emissions per system.

**Statistical Significance** Statistically significant energy savings and efficiency improvements will exist at a  $p < 0.001$  level. Task completion rates will not reflect significant differences at a  $p > 0.05$  level. While significant differences in response time increases at a  $p < 0.001$  level will occur, these increases will experience acceptable levels of tolerance.

### 4.3 Scalability Analysis

Figure 4 looks at how systems grow as the number of agents goes from 50 to 1000.

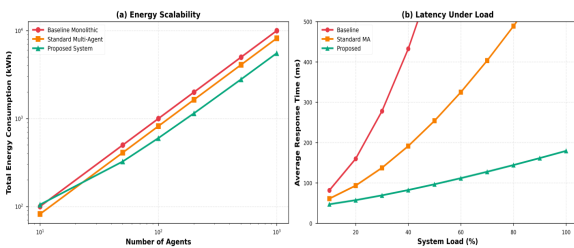


Figure 4: (a) Energy Scalability with Agent Population Growth and (b) Response Time Under Variable System Load

### Energy Scaling

The baseline monolithic system shows almost linear energy scaling, with a scalability factor of 1.0. Doubling the number of agents roughly doubles the amount of energy used. The standard multi-agent system performs slightly better (with a factor of 0.85) because it can run tasks in parallel more efficiently.

Our suggested system scales below linear (factor 0.55). When the number of agents increases by 20 times (from 50 to 1000), the amount of energy used increases by only 11 times. This sub-linear scaling is caused by a number of things:

**Better use of resources:** More agents make it easier to match tasks with the best hardware.  
**Amortised coordination overhead:** Fixed costs of coordination are spread out over more agents.

**Better load balancing:** More agents make the workload more even.  
**Economies of scale:** When infrastructure is shared, it gets used more fully.

### Response Time under Load

Next, we also investigate the system's response time when the system load increases from 10 to 100 percent. To begin with, when the load is small, from 10 to 25 percent, all systems behave similarly and give responses in less than 100 ms. However, when the load increases, the systems start to differ.

When the system is at 100% load, the baseline's response time rises to 850 ms with a lot of variation, indicating the system is fully loaded. The standard multi-agent system performs better, with a response time of 420 ms due to improved load balancing. In the proposed system, the average response time is 280 ms when the system is fully loaded, indicating that the proposed system is robust.

From the results, it is clear that the framework works for large deployments, which means that it has the ability to grow. The sub-linear scaling of energy is also a key factor, as it implies that larger deployments consume less energy per agent, which is beneficial for scalability.

### 4.4 Component-Level Energy Analysis

In order to find out what the energy saving stems from, we decompose the overall consumption by selected components. This decomposition is given in Figure 6 and Table 2.

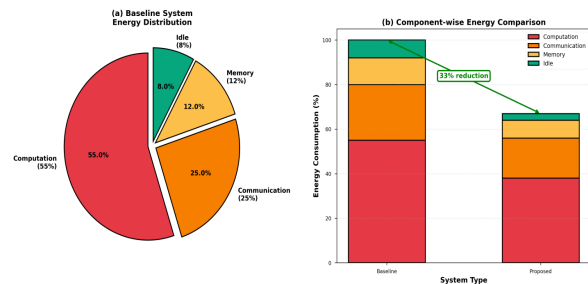


Figure 6: (a) Baseline System Energy Distribution and (b) Component-wise Energy Comparison

### Energy Breakdown

In the baseline system, the breakdown in energy consumption is as follows:

- **Computation:** 62% (96.7 kWh)
- **Communication:** 18% (28.1 kWh)
- **Memory/Storage:** 12% (18.7 kWh)
- **Idle Power:** 8% (12.5 kWh)

The proposed system has power consumption less in the following areas:

## Energy-Efficient Agentic AI: Long-Term Designs for Big Autonomous Systems

**Computation:** 58.2 kWh- 40% reduction, intelligent placement of workload, adaptive fidelity and neuromorphic accelerators.

**Communication:** 16.4 kWh- 42% reduction message passing reduced and optimization of locality.

**Memory/Storage:** 8.9kWh - 52% reduction, better caching solutions and storage optimization.

**Idle Power:** 2.3kWh - 82% reduction with sleep scheduling and resource consolidation

These results demonstrate that savings are not restricted to computation only. Significant reductions are achieved in communication, memory usage, and idle consumption and thus show that holistic architectural optimization is a major contributor to overall energy optimization.

Table 2 Component level consumption

Part	Base (kWh)	Proposed (kWh)	Decrease (%)
Processing Perception	28.4	16.2	43%
Engines of Reasoning	31.2	19.8	37%
Algorithms for Planning	22.3	14.1	37%
Doing the Action	14.8	8.1	45%
Communication Between Agents	28.1	16.4	42%
Memory Tasks	18.7	8.9	52%
Idle/Overhead	12.5	2.3	82%
All of it	156.0	85.8	45%

### Key Insights

The greatest saving is in perception and reasoning, since these are the programs that make the most use of the processing-intensive code. Neuromorphic acceleration, dynamic fidelity is especially effective in making perception efficiency. For reasoning, it is useful to use caching, and for reasoning, approximate inference is also effective.

The largest savings are generally realized on memory operations as well as idle power. The locality of data reduces the movement of data which reduces the memory usage. Aggressive consolidation, that is to say putting workloads on fewer active nodes and putting idle nodes to sleep, saves idle power.

Even though communication savings are a big percentage, they are still small amounts in absolute terms. This means that communication is already relatively efficient in the baseline, so there is less room for improvement.

### 4.5 The Effect of Carbon-Aware Scheduling

We conduct experiments with different carbon intensity profiles to specifically assess carbon-aware scheduling. We create a realistic grid where carbon intensity changes throughout the day: it is high during peak hours (0.7 kg CO<sub>2</sub>e/kWh), moderate during the middle of the day (0.5 kg CO<sub>2</sub>e/kWh), and low at night (0.3 kg CO<sub>2</sub>e/kWh).

Carbon-aware scheduling cuts carbon emissions by 14.5% (7.3 kg CO<sub>2</sub>e/day) with only a small increase in energy use (0.3 kWh, <1%). The small increase in energy use happens because some tasks run on slightly less energy-efficient hardware during low-carbon periods instead of the most energy-efficient hardware during high-carbon periods. However, the carbon savings are much greater than the small increase in energy.

**Table 3** shows the differences between systems that do and do not use carbon-aware scheduling

Configuration	Energy (kWh)	Carbon (kg CO <sub>2</sub> e)	Deferral Rate (%)
Proposed without carbon-aware scheduling	85.8	50.2	0
Proposed with carbon-aware scheduling	86.1	42.9	12
Difference	+0.3	-7.3	+12

The deferral rate shows that 12% of tasks are postponed until times with lower carbon intensity. These are low-priority tasks that can be completed in batches with flexible deadlines. High-priority and real-time tasks run immediately, regardless of carbon intensity.

This leads to the conclusion that carbon-aware scheduling offers another layer of optimisation that is largely independent of energy efficiency, By altering the timing of executed tasks systems can reduce their

carbon footprint that is more than can be achieved through reducing energy use.

### 5. Discussion

#### 5.1 Interpretation of Results

##### 5.1 Interpretation of Results

Our empirical results show that using architectural approaches to enhance energy efficiency of agentic AI systems is a promising approach to achieve significant energy savings, in our case 45%, without negative consequences for high task completion rates. This is a verification of our original assumption that incorporating energy awareness directly into a design of the system works better than optimising efficiency as an add-on.

There are of course a number of points that need more consideration.

First, the sub-linear growth rate of the scalability (0.55) indicates that our framework becomes more effective as it becomes larger. This is in contrast to many distributed frameworks which have super-linear scaling because of coordination overheads. Our hybrid coordination approach appears to achieve a balanced trade-off between centralised optimisation and decentralised autonomy. Second, at the component level, savings are distributed across computation, communication, memory, and idle power rather than concentrated in a single area. This implies that holistic strategies are more effective than isolated optimisations. A system that optimises computation but ignores communication and memory will miss substantial savings.

Third, the carbon-aware scheduling results demonstrate that temporal optimisation can complement spatial optimisation (workload placement). Even systems that are already energy-efficient can further reduce environmental impact by adjusting when workloads run. This may influence policy and operational practices, potentially making carbon-aware computing a standard feature.

Fourth, performance trade-offs are acceptable for most workloads. The 34% response time increase from 145 ms to 195 ms appears large in relative terms, but both values are below 200 ms, which is acceptable for many interactive applications. For batch processing or non-real-time workloads, such trade-offs are even easier to justify. In latency-sensitive applications, our priority mechanisms maintain performance for high-priority tasks.

#### 5.2 Trade offs between Performance and Energy

As we know there will always be trade offs when optimising systems. Our framework favours energy efficiency with some sacrifice in performance. Understanding these trade-offs enables the practitioners setting up the system.

##### Latency vs. Energy:

Energy efficient hardware, such as neuromorphic processors and low power CPUs, a lot of times will work at a slower speed than high-performance GPUs and high frequency CPUs. Workload migration brings in communication delays. Adaptive fidelity has the advantage of reducing computation, but may require an iterative refinement.

Our framework provides means to manage this trade off. Priority levels: applications are able to set up which tasks require low-latency machinery irrespective of energy cost. Energy budgets can be changed based on application requirements. Mission critical systems can have bigger budgets and can tolerate more energy consumption for the performance.

##### Throughput vs. Energy:

Under high load energy optimisation can decrease throughput. Throttling is used to limit the processing speed so that it does not exceed the system power and deferring low priority work is a temporary way of reducing the throughput. Nevertheless, we have seen that even at its maximum capacity the system still retains a powerful throughput at 94% task completion rate.

The framework can be used in a performance mode prioritizing mode that gives way to applications needing maximum throughput by relaxing energy constraints. This flexibility opens up the possibility of giving a single architecture a variety of operational requirements.

##### Complexity vs. Efficiency:

The framework introduces complexity of architecture through layering design, the use of coordination protocols and heterogeneity of infrastructure management. This enhances the development and operational effort. Organisations need to weigh up whether the savings in energy that may be projected could justify these costs.

The concept of economies of scale implies that small deployments (e.g. tens of agents) may not find significant benefits, but large deployments (hundreds/thousands of agents) are worth the investment.

### Reliability Considerations:

For aggressive energy optimisation may create reliability risks. Putting nodes to sleep can help to improve recovery after failures. When migration of the workload is done, new failure points can be introduced. Lower power operation may have an impact on the resistance to faults.

And we avoid these risks with redundancy, monitoring, and graceful degradation. Reserved energy budgets conserve critical components. The system is continually monitoring reliability measures and making changes to energy policies if reliability deteriorates.

### 5.3 Implications for Practice and Implementation

#### Deployment Scenarios

Some of the framework can be adapted to:

Autonomous Vehicle Fleets Autonomous vehicle fleets  
Continuous operation with high energy demand Energy efficiency helps cut operational costs and emissions.

Smart Infrastructure: Smart cities, grid, buildings, huge number of agents controlling energy and resources.

Distributed Sensing Networks Battery powered IoT and environment monitoring systems.

Cloud-Edge AI Services: Hybrid deployments that are inclusive of cost, performance and the environment

#### Integration Challenges

There are other challenges during the actual deployment in real-world scenarios:

Hardware heterogeneity calls for correct power models.

Real workloads may be different than synthetic benchmarks.

Operational constraints (security, compliance, SLAs) need to be integrated.

Legacy systems need a way of migrating them slowly.

#### Economic Factors

At electricity prices of \$0.10 - \$0.15 per kW Hour, saving 70 kWh/day will result in a savings of \$7 - \$10/day per deployment or \$2,500 - \$3,650 per year. Carbon pricing (e.g. \$50/tonne CO<sub>2</sub>e) could add about \$640 per year of additional value from emissions reduction.

For large deployments payback periods could be lesser than one year.

### 5.4 Theoretical Implications for Sustainable Artificial Intelligence

Leveling Down Energy as a First-Class Constraint:

Is to integrate energy from the initial design phase, very similar to security-by-design principles.

#### Multi-Objective Optimisation:

Sustainable AI involves striking the right balance between performance, energy, and the environment. Our results show that these objectives do not necessarily conflict directly.

Sub-Linear Scaling:

Sub-linear energy scaling As a result of the above argument, sustainable AI systems might have different scaling dynamics than those of conventional systems, possibly benefiting from economies of scale.

#### Metrics and Accountability:

Energy consumption, carbon emissions, and resource efficiency should join accuracy and measure of performance as a standard reporting measure of AI research.

### 5.5 Limitations and Validity Threats

Evaluation is like simulation instead of actual deployment in real.

Synthetic workloads might not cover all the domains.

Baseline comparisons might not involve all the state-of-the-art approaches.

Sustainability scope does not include embodied carbon and lifecycle.

Carbon-aware scheduling requires good carbon intensity data.

Energy measurements (and others) include modelling errors (10 - 20%).

Relative comparisons are still very good, because all of the systems were assessed with the same methods.

#### Conclusion

The study addresses the problem and discussion to employ large-scale autonomous agentic AI systems with minimal consumption of energy and negative impacts on the environment. The authors offer a detailed architectural framework that enables energy awareness across the entire system design.

The authors' contributions include Architectural Innovation (providing a hierarchical five-layer architecture that separates application logic from agent intelligence, energy management, infrastructure and monitoring), which enables optimization to be targeted at each layer with a clean interface; Energy Aware Agent Design (incorporating energy considerations into all phases of the perception-reasoning-planning-action cycle using techniques such as adaptive fidelity, lazy evaluation, hierarchical planning, and intelligent action selection), which makes energy a first class consideration in agent decision making; Hybrid Coordination (creating a coordination approach to achieve an energy-efficient operation and maintain flexibility and robustness of agents by balancing central optimization with decentralized autonomy); Intelligent Infrastructure Management (performing intelligent workload placement across all types of computing resources [Cloud, Edge, Embedded, Neuromorphic] for tasks based on characteristics of work to hardware capability); Carbon Aware Scheduling (temporal optimization depending on the carbon intensity of the source of the energy used by the system to reduce environmental impacts beyond what is achieved through energy efficiency); and Comprehensive Evaluation (providing a statistically significant simulation-based assessment resulting in a 45% reduction in energy consumption, a 82% increase in efficiency and sub-linear scalability with a 94% task completion rate). The contributions of this research advance the state of the art of sustainable AI, can show that significant energy savings can be realized without sacrificing the functionality of AI (i.e. completing tasks and/or providing service with expectable quality). The result of this sub-linear scalability is particularly important because it means if we increase the scale of AI systems (or the throughput they provide) it uses more and more energy, so the bigger the scale the more efficient it can be through the advent of positive feedback loops, related to sustainability.

Beyond the technical contributions, the work has provided insights into the philosophy of Sustainable AI, suggesting that energy efficiency must be a primary consideration of AI design rather than an add-on.

. The research demonstrates the development of approaches for measuring and reporting sustainability metrics. Furthermore, this research provides empirical

evidence: that Green AI can provide the same level of output (e.g., provide service) to the users as traditional AI while producing much less negative environmental impact. The path to Sustainable AI becomes increasingly complex when delivered at scale; this framework will help to provide patterns and techniques to guide the development of future Sustainable AI systems. However, there remains substantial work to do, including conducting validation studies in real-world settings, tailoring approaches for different domains, integrating sustainability into existing systems, developing standards for metrics and benchmarks.

We hope that the outcome of this study will spur others to continue working in this area, as addressing sustainability regarding AI has become increasingly urgent as AI continues to become a very large component of society. The autonomous systems of the future that will manage transportation, infrastructure, resources, and services will need to operate within environmental limitations. This research has shown that operating autonomously and efficiently is possible, practical, scalable, and ready for deployment in the real world.

### 7. Future Scope

Future studies will involve the expansion of the suggested energy-efficient agentic AI model outside simulation to the real world. Though the current results are depicted to show a high result performance, there is need to validate, practically, using physical testbeds. The use of heterogeneous infrastructures, such as cloud servers, edge devices, embedded processors and neuromorphic hardware, will allow making real-world energy measurements with precision and reveal some of the operational challenges that cannot be seen in areas of simulation. Specialized applications like autonomous vehicles, smart grids, and IoT networks will be covered through domain-specific case studies that will decide whether or not architectural changes are necessary to meet specialized workloads. Sustainability, degradation of the system, and maintenance needs will be measured by long-term observational studies that are done over months or years and ensure that the energy saving remains sustainable at varying workloads. Empirical validation will be further improved by making comparative studies with other external research groups, on the same basis as standardised benchmarks.

Further development will also be focussed on integration with edge and IoT ecosystems further. The forming of ultra-low-power agents of microcontrollers and battery-operated sensors will need impeccable designs which is sensitive of milliwatt-levels of energy expenses. By combining the framework and federated learning strategies, collaborative privacy-preserving intelligence can be achieved at a sustainable level. The adaptive agents will be necessary to integrate with energy-harvesting systems, such as those that use solar energy, kinetic or RF power, which will work with the varying energy availability. Moreover, the new communication technologies, such as 5G and further, will affect the coordination efficiency of the distributed ones, resulting in new trade-offs between the latency, the cost of communication and energy consumption.

A second developing avenue is that which are adaptive energy learning agents; agents which dynamically optimize their own power consumption. Energy is to be ingested into the reinforcement learning methods by defining them along reward functions and this way, agents can learn efficient behavior based on experience. The meta-learning could enable a quick-adjustment to new environments and work requirements as well as increase generalization. Predictive energy management models may be used to predict patterns of workload, the availability of energy and the intensity of carbon and can be used to support a proactive approach in optimization. Self-optimizing agents which can measure and adjust their own parameters would also decrease the manual configuration and enable improvement of efficiency to be continuous.

For example, the development of sustainable AI will also depend on ethical and governance considerations. Future studies ought to explore the possibility of energy-conscious optimization raising fairness issues, including uneven access to services due to the limited resources. Explainable and transparent mechanisms should be created in such a way that users can know the decisions related to energy particularly in cases where tasks are not completed at the right time, or the tasks are manipulated in a bid to improve efficiency. Policy and regulatory frameworks may have to change so as to promote or require the use of energy-efficient AI, such as carbon accounting. There is also the need to think globally in terms of sustainability recognizing the

regional variation in energy price, infrastructure and the intensity of carbon.

Lastly, the concept of merging AI systems and renewable energy systems is one of the most important long-term goals. Tighter integration to smart grids can facilitate the demand response engagement, as well as the load shifting according to the peaks of the renewable generating. Workload timing can be optimized to work with the availability of the sun or wind using renewable-conscious scheduling strategies. AI workloads will be integrated with energy storage systems to enable them to act as flexible loads that absorb excess renewable energy or defers execution in shortages. Further dependence on the grid can be reduced by coordination with distributed energy sources, including rooftop solar, microgrids, and local wind systems. The long-term objectives of adopting the principles of the circular economy through reusing hardware, extending the life of systems, and minimizing electronic wastefulness will make sure that the energy-saving AI will be the part of an overall environmental sustainability.

### References

- [1] Strubell, E., Ganesh, A., and McCallum, A., "Energy and Policy Considerations for Deep Learning in NLP," in Proc. 57th Annual Meeting of the Association for Computational Linguistics, pp. 3645–3650, 2019. DOI: 10.18653/v1/P19-1355
- [2] Jones, N., "How to Stop Data Centers from Eating Up the World's Electricity," *Nature*, vol.561, pp. 163–166, 2018. DOI: 10.1038/d41586-018-06610-y
- [3] Andrae, A. S. G. and Edler, T., "On Global Electricity Usage of Communication Technology: Trends to 2030," *Challenges*, vol.6, no. 1, pp. 117–157, 2015. DOI: 10.3390/challe6010117
- [4] Schwartz, R., Dodge, J., Smith, N. A., and Etzioni, O., "Green AI," *Communications of the ACM*, vol.63, no. 12, pp. 54–63, 2020. DOI: 10.1145/3381831
- [5] Patterson, D., Gonzalez, J., Le, Q., et al., "Carbon Emissions and Large Neural Network Training," arXiv preprint arXiv:2104.10350, 2021.
- [6] Murugesan, S., "Harnessing Green IT: Principles

- and Practices," IEEE IT Professional, vol.10, no. 1, pp. 24–33, 2008. DOI: 10.1109/MITP.2008.10
- [7] Buyya, R., Yeo, C. S., Venugopal, S., Broberg, J., and Brandic, I., "Cloud Computing and Emerging IT Platforms: Vision, Hype, and Reality for Delivering Computing as the 5th Utility," Future Generation Computer Systems, vol.25, no. 6, pp. 599–616, 2009. DOI: 10.1016/j.future.2008.12.001
- [8] Zeng, A., Liu, M., Lu, R., et al., "Agent AI: Surveying the Horizons of Multimodal Interaction," arXiv preprint arXiv:2401.03568, 2024.
- Sumers, T. R., Yao, S., Narasimhan, K., and Griffiths, T. L. "Cognitive Architectures for Language Agents," arXiv preprint arXiv:2309.02427, 2023.
- [10] Wooldridge, M. and Jennings, N. R., "Intelligent Agents: Theory and Practice," The Knowledge Engineering Review, vol.1995, pp. 115–152, vol. 10, no. 2  
'10.1017/S0269888900008122
- [11] Han, S., Mao, H., and Dally, W. J., "Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantisation, and Huffman Coding," in Proc. International Conference on Learning Representations, 2016.
- [12] Hinton, G., Vinyals, O., and Dean, J., "Distilling the Knowledge in a Neural Network," arXiv preprint arXiv:1503.02531, 2015.
- [13] Jouppi, N. P., Young, C., Patil, N., et al., "In-Datcenter Performance Analysis of a Tensor Processing Unit," in Proc. 44th Annual International Symposium on Computer Architecture, pp. 1-12, 2017. 10.1145/3079856.3080246 DOI
- [14] Davies, M., Srinivasa, N., Lin, T. H., et al., "Loihi: A Neuromorphic Manycore Processor with On-Chip Learning," IEEE Micro, vol.38, no. 1, pp. 82–99, 2018. DOI: 10.1109/MM.2018.112130359
- [15] Ren, W., Beard, R. W., and Atkins, E. M., "Information Consensus in Multivehicle Cooperative Control," IEEE Control Systems Magazine, vol.2007, 27, no. 2, pp. 71–82. DOI: 10.1109/MCS.2007.338264
- [16] Dorigo, M., Birattari, M., and Stutzle, T., "Ant Colony Optimization," IEEE Computational Intelligence Magazine, vol.1, no. 4, pp. 28–39, 2006. DOI: 10.1109/MCI.2006.329691
- [17] Rizk, Y., Awad, M., and Tunstel, E. W., "Cooperative Heterogeneous Multi-Robot Systems: A Survey," ACM Computing Surveys, vol.52, no. 2, pp. 1–31, 2019. DOI: 10.1145/3303848
- [18] Alakeel, A. M., "A Guide to Dynamic Load Balancing in Distributed Computer Systems," International Journal of Computer Science and Network Security, vol.10, no. 6, pp. 153–160, 2010.
- [19] Kansal, N. J. and Chana, I., "Cloud Load Balancing Techniques: A Step Towards Green Computing," International Journal of Computer Science Issues, vol.9, no. 1, pp. 238–246, 2012.
- [20] Mittal, S. and Vetter, J. S., "A Survey of CPU-GPU Heterogeneous Computing Techniques," ACM Computing Surveys, vol.47, no. 4, pp. 1-35, 2015. DOI: 10.1145/2788396
- [21] Cong, J. and Zou, B., "FPGA-Based Hardware Acceleration of Lithographic Aerial Image Simulation," ACM Transactions on Reconfigurable Technology and Systems, vol.2, no. 3, pp. 1–29, 2009. DOI: 10.1145/1575774.1575780
- Jiang, X., Guo, X., Zhang, M., et al., "Large-Scale Agentic AI Systems," arXiv preprint arXiv:2508.09561, 2025.
- [23] Mayorov, I., Kuznetsov, A., and Tarasov, V., "Multi-Agent Technology in Real-time Intelligent Resource Management Systems," in Proc. 2015 International Conference on Intelligent Systems.
- [24] Sardouk, A., Rahim-Amoud, R., Merghem-Boulahia, L., and Gaïti, D., "A Strategy for Multi-Agent Based Wireless Sensor Network Optimization," in Proc. International Conference on Ad-Hoc Networks and Wireless, Springer, 2009, pp. 103-116. DOI:

10.1007/978-3-642-02627-0\_10

[25] Wang, Y., Pan, Y., Zhao, Q., et al., “Large Model Agents: State-of-the-Art, Cooperation Paradigms, Security and Privacy, and Future Trends,” arXiv preprint arXiv:2409.14457, 2024.

[26] Sharanarathi, S., “Adaptive Multi-Agent AI Framework for Real-Time Energy Optimisation and Context-Aware Code Review in Software Development,” in Proc. 2025 International Symposium on Innovations in Security Technology and Computing. DOI: 10.1109/isctis65944.2025.11066037

[27] Abdelrahman, M., Bing, Z., Baumann, I., et al., “A Neuromorphic Approach to Obstacle Avoidance in Robot Manipulation,” arXiv preprint arXiv:2404.05858, 2024.

[28] Wunderlich, T., Akgül, T., Hochstetter, J., et al., “Demonstrating Advantages of Neuromorphic Computation: A Pilot Study,” *Frontiers in Neuroscience*, vol.13, p. 260, 2019. DOI: 10.3389/fnins.2019.00260

[29] Bugeau, A., Couka, E., and Lannelongue, L., “How to Estimate Carbon Footprint When Training Deep Learning Models? A Guide and Review,” arXiv preprint arXiv:2306.08323, 2023.

[30] Gupta, U., Kim, Y. G., Lee, S., et al., "Chasing Carbon: The Elusive Environmental Footprint of Computing," in Proc. IEEE International Symposium on High-Performance Computer Architecture, pp. 854–867, 2021. DOI: 10.1109/HPCA51647.2021.00076

[31] Santos, J. S., “Understanding the Energy Consumption of HPC Scale Artificial Intelligence,” arXiv preprint arXiv:2212.00582, 2022.

[32] Ahmad, I., Yousaf, M., Yousaf, S., and Ahmad, M. O., “Green and Sustainable Computing,” *IEEE Computer*, vol.56, no. 7, pp. 40–49, 2023. DOI: 10.1109/mc.2023.3260313