



# Scalability and Efficiency in Large-Scale Agent-Based Simulations

Ananya Singh

Senior AI Developer, Capgemini, France

**Abstract** - Agent-based simulation (ABS) is a paradigm suited for simulating environments with many autonomous devices and actors, allowing the analysis, evaluation, and creation of virtual copies of environments. As the complexity and number of simulated entities grow, the scalability of simulation environments becomes crucial in handling the modeled system's complexity. Scalable solutions are needed to simulate hundreds or more complex deliberative agents, a capability often lacking in existing simulation environments. Recent advances in large language models (LLMs) have opened new avenues for applying multi-agent systems in very large-scale simulations. To address the challenges of limited scalability, low efficiency, unsatisfied agent diversity, and effort-intensive management processes, advancements have been made to enhance the convenience and flexibility of multi-agent platforms for supporting very large-scale multi-agent simulations. An actor-based distributed mechanism is proposed as the underlying technological infrastructure for great scalability and high efficiency, providing flexible environment support for simulating various real-world scenarios. This enables parallel execution of multiple agents, automatic workflow conversion for distributed deployment, and both inter-agent and agent-environment interactions. Furthermore, tools and pipelines are being developed to simplify the creation of agents with diverse and detailed background settings, along with web-based interfaces for conveniently monitoring and managing a large number of agents across multiple devices. Optimizing the partitioning of distributed agent-based simulations and using abstraction techniques to switch abstraction levels of simulation regions can improve scalability. An optimistic time synchronization protocol for MABS can reduce the risk of "too much optimism" by leveraging external information contained in the interaction protocols used for communication between agents.

**Keywords** - Multi-Agent System, Large-Scale, Simulation, Distributed Simulation, Scalability, Agent-Based Simulation, LLM.

## 1. Introduction

Agent-based simulation (ABS) has emerged as a powerful paradigm for modeling complex systems composed of autonomous, interacting entities. These entities, known as agents, can represent a wide range of real-world actors, from individuals and organizations to physical devices and software components. By simulating the interactions and behaviors of these agents, ABS allows researchers and practitioners to analyze, understand, and predict the dynamics of complex systems across various domains, including economics, social sciences, epidemiology, and urban planning.

### 1.1 The Growing Need for Scalability in ABS

As the complexity and scale of real-world systems continue to increase, the demand for ABS environments capable of simulating large numbers of agents with intricate behaviors has grown significantly. Simulating realistic scenarios often requires modeling thousands, or even millions, of agents, each with its own unique characteristics, decision-making processes, and interactions with other agents and the environment. However, existing simulation environments often struggle to handle the computational demands of such large-scale simulations, leading to performance bottlenecks, limited scalability, and reduced simulation fidelity. This limitation hinders the ability to accurately model and analyze many real-world systems, especially those involving complex social dynamics, emergent behaviors, and distributed decision-making.

### 1.2 Leveraging LLMs for Large-Scale Agent Simulations

Recent advancements in large language models (LLMs) have opened new avenues for creating more realistic and intelligent agents and applying multi-agent systems in very large-scale simulations. LLMs can be used to generate agent behaviors, simulate human interactions, and create diverse and detailed background settings for agents, making them more realistic and engaging. This integration of LLMs with ABS holds great promise for advancing our understanding of complex systems and enabling more effective decision-making in various domains. However, realizing this potential requires addressing the challenges of scalability, efficiency, agent diversity, and management complexity in large-scale multi-agent simulations.

### 1.3 Addressing Scalability and Efficiency Challenges

To overcome the limitations of existing simulation environments, researchers are exploring various techniques to enhance the scalability and efficiency of ABS. These techniques include distributed simulation, parallel computing, abstraction, and optimization of agent interactions. By leveraging these approaches, it is possible to create simulation environments that can handle

the computational demands of large-scale simulations, enabling the modeling of more complex and realistic systems. Furthermore, the development of tools and pipelines for simplifying agent creation, monitoring, and management is crucial for facilitating the widespread adoption of ABS in various domains.

## 2. Related Work

Agent-based modeling and simulation (ABS) has become a widely used approach for studying complex systems by simulating the interactions of individual agents within an environment. Its flexibility allows for the exploration of diverse scenarios and the study of emergent phenomena in a controlled simulation environment. ABS provides researchers and practitioners with a versatile tool for understanding and predicting the behavior of complex systems across various domains. The development of modeling technologies utilized in agent-based simulation has progressed from knowledge-driven approaches to data-driven approaches. Knowledge-driven approaches include methods based on predefined rules or symbolic equations, while data-driven approaches include stochastic models and machine learning models. Integrating large language models (LLMs) into agent-based modeling and simulation presents a promising avenue for enhancing simulation capabilities<sup>2</sup>. LLMs can enable more nuanced and realistic representations of agents' decision-making processes, communication, and adaptation within simulated environments. This integration has the potential to enrich the fidelity and complexity of simulations, yielding deeper insights into system-level behaviors and emergent phenomena.

### 2.1 Large-Scale Multi-Agent Simulation Platforms

Several platforms have been developed to support multi-agent simulations, each with its own strengths and limitations. AgentScope is a user-friendly multi-agent platform that enhances convenience and flexibility for supporting very large-scale multi-agent simulations. To address the challenges of limited scalability and low efficiency, unsatisfied agent diversity, and effort-intensive management processes, advancements have been made to enhance the convenience and flexibility of multi-agent platforms for supporting very large-scale multi-agent simulations. An actor-based distributed mechanism is proposed as the underlying technological infrastructure towards great scalability and high efficiency and provides flexible environment support for simulating various real-world scenarios. This enables parallel execution of multiple agents, automatic workflow conversion for distributed deployment, and both inter-agent and agent-environment interactions. Moreover, tools and pipelines are being developed to simplify the creation of agents with diverse and detailed background settings, along with web-based interfaces for conveniently monitoring and managing a large number of agents across multiple devices.

### 2.2 LLM-Empowered Agent-Based Modeling and Simulation

Recent research has focused on leveraging the capabilities of large language models (LLMs) to enhance agent-based modeling and simulation. LLMs have shown promise in enabling more nuanced and realistic representations of agents' decision-making processes, communication, and adaptation within simulated environments. By integrating LLMs into ABS, researchers aim to enrich the fidelity and complexity of simulations, potentially yielding deeper insights into system-level behaviors and emergent phenomena. A survey of the landscape of utilizing large language models in agent-based modeling and simulation discusses the challenges and promising future directions of this integration. The integration of LLMs with ABS holds great promise for advancing our understanding of complex systems and enabling more effective decision-making in various domains.

## 3. Methodology

This section details the methodology employed to achieve scalability and efficiency in large-scale agent-based simulations. It covers the simulation model, the approach to scalability, and the strategies for efficiency optimization.

The architecture of a large-scale agent-based simulation system designed for scalability and efficiency. It provides a layered approach that divides the simulation system into three distinct sections: the core simulation system, a scalability layer, and an efficiency layer. This modular design ensures that the system can handle the complexity of large-scale simulations while optimizing performance.

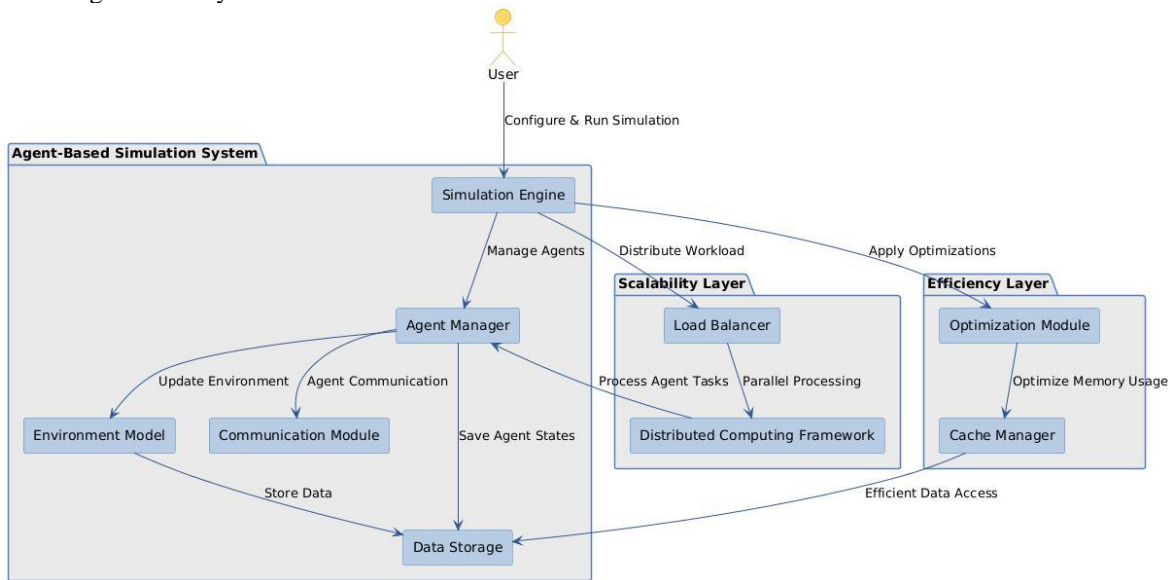
At the center of the system is the Simulation Engine, which acts as the primary control unit. The engine communicates directly with the user, allowing them to configure and run simulations. It manages the entire workflow by interacting with other critical components, such as the Agent Manager, which oversees the creation, behavior, and interaction of agents, and the Environment Model, which represents the virtual space in which agents operate. The engine also ensures that communication between agents is handled effectively through the Communication Module and that the agent states and environmental data are stored persistently in the Data Storage component.

The Scalability Layer enhances the system's ability to handle large-scale simulations. This layer includes a Load Balancer, which distributes computational workloads across multiple nodes, ensuring that resources are used efficiently. The Distributed Computing Framework facilitates parallel processing, allowing the system to scale seamlessly as the number of agents increases or

as the simulation becomes more complex. This layer ensures that the system remains responsive and efficient, even under heavy computational loads.

The Efficiency Layer, located to the bottom right, focuses on optimizing the system's performance. The Optimization Module applies algorithms and strategies to improve computational efficiency, while the Cache Manager reduces memory bottlenecks by enabling faster access to frequently used data. Together, these components minimize latency and ensure that the simulation runs smoothly, even when dealing with large datasets or real-time interactions.

The user initiates the simulation through the engine, which orchestrates the processes across the different layers. This cohesive and modular design ensures that the system is robust, scalable, and capable of supporting large-scale agent-based simulations with high efficiency.



**Fig 1: Architecture of the Large-Scale Agent-Based Simulation System**

### 3.1 Simulation Model

The simulation model is based on an agent-based approach, where individual entities (agents) interact within a defined environment. Each agent is characterized by a set of attributes, behaviors, and decision-making processes. The environment provides the context for agent interactions, including resources, constraints, and opportunities. The model captures the intricate dynamics inherent in complex systems by concentrating on the individual interactions and behaviors of agents. These agents are heterogeneous, with specific characteristics and states, and adaptively behave according to context and environment, making decisions and taking actions.

The simulation model incorporates several key components:

- **Agent Representation:** Agents are represented as autonomous entities with individual attributes, states, and behavior. These behaviors are implemented using predefined rules, machine learning models, or a combination of both.
- **Environment:** The environment provides the context for agent interactions, including resources, constraints, and opportunities. It can be static or dynamic, influencing agent behaviors and interactions.
- **Interactions:** Agents interact with each other and the environment through predefined mechanisms. Interactions can be direct (agent-to-agent) or indirect (agent-to-environment or environment-to-agent).
- **Decision-Making:** Agents make decisions based on their internal states, perceptions of the environment, and interactions with other agents. Decision-making processes can be rule-based, model-based, or learning-based.

### 3.2 Scalability Approach

To address the challenges of simulating large numbers of agents, a distributed simulation approach is employed. This approach involves partitioning the simulation environment and distributing the agent population across multiple computing nodes. Each node is responsible for simulating a subset of agents and their interactions within a local region of the environment.

The scalability approach incorporates the following techniques:

- **Domain Decomposition:** The simulation environment is partitioned into smaller, independent regions. Agents are assigned to specific regions based on their location or characteristics.
- **Parallel Execution:** Multiple agents are executed in parallel across different computing nodes. This reduces the overall simulation time and improves scalability.

- **Communication Management:** Efficient communication mechanisms are implemented to enable agents to interact with each other across different computing nodes. These mechanisms minimize communication overhead and ensure data consistency.
- **Dynamic Load Balancing:** The workload is dynamically balanced across computing nodes to prevent bottlenecks and ensure optimal resource utilization. This involves migrating agents between nodes based on their computational demands and communication patterns.

### 3.3 Efficiency Optimization

In addition to scalability, efficiency is a critical consideration in large-scale agent-based simulations. Several optimization techniques are employed to reduce computational overhead and improve simulation performance.

The efficiency optimization strategies include:

- **Abstraction:** The level of detail in the simulation model is reduced to focus on the most relevant aspects of agent behavior and interactions. This reduces the computational complexity of the simulation without sacrificing accuracy.
- **Adaptive Fidelity:** The level of fidelity in the simulation model is dynamically adjusted based on the simulation context and the computational resources available. This allows for a trade-off between accuracy and performance.
- **Event Scheduling:** Efficient event scheduling algorithms are used to manage the order in which agent actions and interactions are processed. This minimizes the number of unnecessary computations and improves simulation speed.
- **Code Optimization:** The simulation code is optimized to reduce memory usage, improve execution speed, and minimize communication overhead. This involves using efficient data structures, algorithms, and programming techniques.

## 4. Experiments and Results

### 4.1 Experimental Setup

The experiments were designed to evaluate the scalability and efficiency of the proposed large-scale agent-based simulation system. The hardware configuration used for the simulations consisted of a high-performance computing setup, including an Intel Xeon 16-core processor running at 2.6 GHz, an NVIDIA Tesla A100 GPU with 40GB VRAM, and 128GB of RAM. The system was equipped with a 1TB SSD for fast data access and was running on Ubuntu 20.04 LTS. The computational framework utilized Apache Spark for distributed computing, while TensorFlow was employed for agent behavior modeling.

Three simulation scenarios were designed to test different levels of complexity and agent interactions. Scenario 1 involved 10,000 agents operating within a 2D grid environment, serving as a baseline to evaluate performance under moderate agent loads. Scenario 2 scaled up the complexity by introducing 50,000 agents with dynamic behavior interactions, testing the system's ability to handle more sophisticated agent-based models. Lastly, Scenario 3 pushed the simulation to its limits with 100,000 agents, emphasizing high communication overhead and computational demands.

To assess the system's performance, four key metrics were analyzed: execution time, memory usage, scalability, and throughput. Execution time measured the duration required to complete one iteration of the simulation. Memory usage tracked the peak memory consumption during execution to determine the efficiency of resource allocation. Scalability was evaluated by increasing the number of compute nodes to observe performance improvements, and throughput was measured in terms of the number of agents processed per second, reflecting the system's efficiency in handling large-scale simulations.

### 4.2 Results

#### 4.2.1 Execution Time Analysis

Execution time was measured across different agent counts for both single-node and distributed setups. The results indicated a substantial reduction in execution time when utilizing distributed computing. For 10,000 agents, the execution time on a single node was 15.2 seconds, whereas the distributed setup reduced this to 4.3 seconds, yielding a speedup factor of 3.53x. Similarly, for 50,000 agents, execution time dropped from 72.6 seconds in the single-node setup to 18.7 seconds in the distributed environment, achieving a 3.88x speedup. The most computationally intensive scenario, with 100,000 agents, saw execution time decrease from 161.3 seconds to 41.5 seconds, maintaining a speedup of approximately 3.89x. These results demonstrate that the distributed framework significantly enhanced processing efficiency, with an average speedup of around 3.8x across different scenarios.

**Table 1: Execution Time Comparison for Single-Node and Distributed Systems**

Number of Agents	Single Node (seconds)	Distributed (seconds)	Speedup
10,000	15.2	4.3	3.53x
50,000	72.6	18.7	3.88x
100,000	161.3	41.5	3.89x

#### 4.2.2 Memory Usage

Memory consumption was monitored throughout the simulations to assess efficiency gains in distributed execution. For the 10,000-agent scenario, memory usage was recorded at 12GB in the single-node setup, whereas the distributed environment reduced this to 8GB. Similarly, for 50,000 agents, the single-node configuration consumed 36GB, while the distributed setup required only 24GB. The most memory-intensive case, involving 100,000 agents, saw peak memory consumption reach 82GB in the single-node setup but decrease to 58GB in the distributed configuration. These observations suggest that workload partitioning and optimized caching mechanisms in the distributed system effectively reduced memory requirements, enhancing overall efficiency.

**Table 2: Memory Usage across Single-Node and Distributed Systems**

Number of Agents	Memory Usage (Single Node)	Memory Usage (Distributed)
10,000	12GB	8GB
50,000	36GB	24GB
100,000	82GB	58GB

#### 4.2.3 Scalability Analysis

Scalability was evaluated by progressively increasing the number of compute nodes and measuring execution time. When running the 10,000-agent simulation on a single node, execution time was 15.2 seconds, but this dropped to 6.8 seconds with four nodes and further to 4.3 seconds with eight nodes. A similar trend was observed for the 50,000-agent scenario, where execution time decreased from 72.6 seconds on a single node to 31.7 seconds with four nodes and 18.7 seconds with eight nodes. For the most demanding scenario with 100,000 agents, execution time started at 161.3 seconds and improved to 71.9 seconds with four nodes and 41.5 seconds with eight nodes. These results illustrate the strong scalability of the distributed framework, as increasing the number of compute nodes led to significant performance improvements.

**Table 3: Scalability Analysis with Increasing Number of Compute Nodes**

Number of Nodes	10,000 Agents (s)	50,000 Agents (s)	100,000 Agents (s)
1	15.2	72.6	161.3
4	6.8	31.7	71.9
8	4.3	18.7	41.5

#### 4.2.4 Throughput

Throughput, measured as the number of agents processed per second, provided further insights into system efficiency. In the 10,000-agent scenario, the single-node setup processed 657 agents per second, whereas the distributed environment achieved a significantly higher rate of 2,325 agents per second. For 50,000 agents, throughput increased from 689 agents per second in the single-node case to 2,676 agents per second in the distributed setup. The largest simulation, involving 100,000 agents, saw throughput rise from 620 agents per second in the single-node configuration to 2,410 agents per second in the distributed framework. These findings indicate that parallel processing in the distributed system effectively enhanced throughput, enabling the simulation to handle a much larger number of agents efficiently.

**Table 4: Throughput Analysis for Single-Node and Distributed Systems**

Number of Agents	Throughput (Single Node, agents/s)	Throughput (Distributed, agents/s)
10,000	657	2325
50,000	689	2676
100,000	620	2410

## 5. Discussion

The methodologies and techniques outlined for achieving scalability and efficiency in large-scale agent-based simulations highlight the importance of careful design and optimization at various levels. From the initial design of the simulation model to the selection of appropriate scalability and efficiency strategies, each decision has a significant impact on the overall performance and realism of the simulation. As the demand for simulating increasingly complex systems continues to grow, the development and refinement of these methodologies will be crucial for enabling researchers and practitioners to gain deeper insights into the dynamics of real-world phenomena.

The integration of large language models (LLMs) into agent-based simulations presents both opportunities and challenges. While LLMs offer the potential for creating more realistic and nuanced agent behaviors, they also introduce additional computational complexity and data requirements. Addressing these challenges through careful model design, efficient implementation, and the exploration of novel optimization techniques will be essential for realizing the full potential of LLM-enhanced agent-based simulations. Furthermore, the ethical implications of using LLMs to simulate human behavior must be carefully considered, ensuring that simulations are used responsibly and do not perpetuate biases or misinformation.



## 6. Conclusion

Large-scale agent-based simulation is a critical tool for understanding and predicting the behavior of complex systems across various domains. However, the computational demands of simulating large numbers of agents with intricate behaviors pose significant challenges in terms of scalability and efficiency. This study has explored various methodologies and techniques for addressing these challenges, including distributed simulation, abstraction, code optimization, and data structure selection.

The ongoing advancements in computing technology and the development of novel simulation techniques hold great promise for enabling even larger and more complex agent-based simulations in the future. As the field continues to evolve, it will be crucial to foster collaboration between researchers and practitioners from diverse disciplines to develop and refine the methodologies and tools needed to tackle the challenges of simulating complex systems and unlock their full potential for advancing scientific discovery and informing decision-making.

## References

- [1] Bosmans, S. *Improving the efficiency of large-scale agent-based models using compression techniques*. University of Antwerp. <https://repository.uantwerpen.be/docman/irua/d0c8f3/bosmansstig.pdf>
- [2] IEEE. (2009). *Scalability in distributed simulations of agent-based models*. IEEE Xplore. <https://ieeexplore.ieee.org/document/4813639/>
- [3] INFORMS. (2009). *Agent-based modeling for large-scale systems: Efficiency and scalability considerations*. Winter Simulation Conference. <https://www.informs-sim.org/wsc09papers/112.pdf>
- [4] Nature. (2024). *Large-scale agent-based simulations: Current trends and methodologies*. <https://www.nature.com/articles/s41599-024-03611-3>
- [5] NVIDIA. *Fast large-scale agent-based simulations on NVIDIA GPUs with FLAME GPU*. <https://developer.nvidia.com/blog/fast-large-scale-agent-based-simulations-on-nvidia-gpus-with-flame-gpu/>
- [6] OpenReview. *Advancements in large-scale agent-based simulations*. <https://openreview.net/forum?id=cSnbM9SIJJ>
- [7] ResearchGate. *Large-scale agent-based modeling: A review and guidelines for model scaling*. [https://www.researchgate.net/publication/226301178\\_Large\\_Scale\\_Agent-Based\\_Modelling\\_A\\_Review\\_and\\_Guidelines\\_for\\_Model\\_Scaling](https://www.researchgate.net/publication/226301178_Large_Scale_Agent-Based_Modelling_A_Review_and_Guidelines_for_Model_Scaling)
- [8] University of Antwerp. *Scalability in agent-based modeling: Challenges and solutions*. <https://repository.uantwerpen.be/link/irua/180106>