

# Foundations of Machine Learning and Breakthroughs in Transformers: The Evolutionary Path Behind AI Agents

Junyu Liu

Department of Computer Science, Macau University of Science and Technology, Macau, China

## Abstract

This article explores the evolutionary trajectory of AI Agents, highlighting the foundational role of machine learning and the transformative impact of the Transformer architecture. It first elaborates on how various machine learning paradigms—including supervised, unsupervised, and reinforcement learning—construct the cognitive, perceptual, and execution frameworks for intelligent agents. Subsequently, the study analyzes the core innovations of the Transformer model, such as self-attention mechanisms and parallel computing, demonstrating its function as the central reasoning engine for handling multimodal data and complex tasks. Furthermore, the paper traces the developmental stages of AI Agents from early rule-based systems to the current era of large model-driven and multi-agent collaborative networks. Finally, it examines pressing technical challenges, such as model interpretability and data bias, while outlining significant application opportunities in fields like intelligent manufacturing, healthcare, and transportation, providing a comprehensive theoretical perspective on the future of autonomous agents.

**Key words:** AI Agent; Technical Challenges; Application Opportunities; Large Models.

## 1. Introduction

### 1.1 The Essence and Capabilities of Machine Learning

Machine learning is a technology that enables computers to automatically learn from data and make predictions or decisions through algorithms and models [1]. Its core objective is to construct models that are adaptable to new data, providing a foundational cognitive framework for AI Agents. Data-driven is an important feature of machine learning, and AI Agents rely on machine learning models to process multimodal data, including text, images, and speech.

Supervised learning is one of the most commonly used methods in machine learning. It trains models using labeled training data, enabling the models to make accurate predictions on new input data [2]. For example, in the smart home scenario, AI Agents use supervised learning to recognize user voice commands and invoke corresponding devices. By training the model with a large amount of labeled voice data, the model can learn the mapping relationship between different voice commands and device operations, thereby achieving accurate recognition and response to user instructions.

Unsupervised learning discovers patterns and structures in data without labels [3]. In the field of financial risk control, unsupervised learning can detect abnormal transaction behaviors. By performing cluster analysis on a large amount of transaction data, the model can classify normal transactions and abnormal transactions into different clusters, thereby identifying potential financial risks.

Reinforcement learning is another important branch of machine learning. It enables agents to learn optimal strategies through interactions with the environment [4]. In reinforcement learning, agents take actions based on the state of the environment and receive rewards or penalties based on the results of their actions. Through continuous trial and error, they optimize their strategies. For example, in autonomous driving, AI Agents use reinforcement learning to adjust decision-making strategies in complex road conditions, improving safety. In simulation environments, agents continuously try different driving strategies, receiving corresponding rewards or penalties based on traffic conditions and safety indicators, and gradually learning the optimal driving strategies in various road conditions.

The adaptive and generalization capabilities of machine learning models enable AI Agents to adapt to dynamic environments. Adaptive capability refers to the model's ability to continuously adjust its parameters based on new data to better adapt to environmental changes. Generalization capability is the ability of the model to perform well on unseen data [5]. For example, in recommendation systems, AI Agents use machine learning models to analyze users' historical behavior data, learning their interest preferences. When new products or content appear, the model can make personalized recommendations based on users' interest preferences, even if these new products or content have not appeared in the training data before.

### 1.2 The Core Role of Machine Learning in AI Agents

Machine learning models play multiple roles in AI Agents. At the perception layer, models such as convolutional neural networks (CNN) and recurrent neural networks (RNN) convert raw data into structured information, providing the AI Agent with the ability to perceive the environment [6]. For example, the Vision Transformer (ViT) captures global features of images through self-attention mechanisms, enabling the AI Agent to recognize objects or scenes [7]. ViT divides images

into multiple small blocks and then inputs these blocks as sequences into the Transformer model, learning the interrelationships between the blocks through self-attention mechanisms to capture the global features of the image. Compared to traditional CNN models, ViT has better performance and generalization capabilities when processing large-scale image data.

At the decision-making layer, reinforcement learning algorithms enable AI Agents to learn optimal strategies in interactions with the environment. For example, game AI optimizes action selection through reinforcement learning to improve winning rates [8]. In the Deep Q-Network (DQN), the agent estimates the Q-values of each action using a neural network, which is the expected value of long-term rewards obtained by taking that action. By continuously updating the parameters of the neural network, the agent can learn the optimal action to take in each state, thereby achieving better results in the game.

At the execution layer, machine learning models access external APIs through function calls or plugin mechanisms to implement specific task execution. For example, intelligent customer service AI Agents call knowledge base APIs to answer user questions [9]. The intelligent customer service system understands user questions through natural language processing models, and then calls corresponding knowledge base APIs based on the type and keywords of the question to obtain relevant answers and return them to the user. This intelligent customer service system based on machine learning models can quickly and accurately answer user questions, improving customer satisfaction.

## 2. Transformer Breakthrough: The Reasoning Engine of AI Agent

### 2.1 Core Innovation of Transformer

As a deep learning architecture based on self-attention mechanism, Transformer was proposed by the Google team in 2017. Since then, it has sparked a revolution in the field of artificial intelligence with its unique innovations. The core innovations of Transformer mainly lie in the self-attention mechanism, multi-head attention, and parallel computing.

The self-attention mechanism is the cornerstone of Transformer. It breaks the limitation of traditional sequence processing models that can only process positions one by one, allowing the model to simultaneously consider information from all positions when processing sequence data. This mechanism enables the model to easily capture long-range dependencies, such as in text generation tasks, where the model can accurately understand the referential relationships of pronouns in the context, thereby generating more coherent and accurate text.

Multi-head attention is an expansion and deepening of the self-attention mechanism. It extends the self-attention mechanism into multiple parallel attention heads, each of which independently learns different attention weights. This design enables the model to capture complex relationships in the sequence from multiple perspectives and levels. For example, in the machine translation task, multi-head attention can simultaneously focus on the grammatical structure and semantic features of both the source and target languages, thereby significantly improving the accuracy and fluency of translation.

Parallel computing is another highlight of Transformer. It supports parallel processing of all positions in the input sequence for the self-attention mechanism, significantly improving the training and inference efficiency of the model. Compared with traditional recurrent neural network models, Transformer not only processes long sequences faster but also has lower memory usage, enabling it to easily handle large-scale data processing requirements.

### 2.2 Application of Transformer in AI Agent

The large language model of the Transformer architecture has become the core brain of AI Agent, supporting a series of complex tasks such as natural language understanding, dialogue generation, and text summarization. For example, the intelligent writing assistant AI Agent can generate article outlines or paragraphs by leveraging the large language model, providing powerful assistance to writers.

In addition to natural language processing tasks, Transformer also supports the joint processing of multimodal data such as text, images, and speech through extended architectures. This multimodal processing capability has greatly enhanced the environmental perception ability of AI Agent. For example, in the intelligent medical AI Agent, it can combine medical images and medical record texts for disease diagnosis, providing doctors with more comprehensive and accurate diagnostic basis.

Furthermore, the combination of Transformer with vector databases provides strong support for the long-term memory management of AI Agent. By retrieving historical conversation records, intelligent assistant AI Agent can provide users with more personalized services, thereby enhancing user satisfaction and loyalty. This long-term memory management capability enables AI Agent to handle complex tasks more proficiently and provide users with a more convenient and efficient usage experience.

## 3. Evolution of AI Agents: From Rule-Based to Autonomous Agents

### 3.1 Early Technological Foundation Period

In the early stage, AI Agents were based on symbolic methods and implemented task execution through preset rules and finite state machines. For example, chatbots based on keyword matching matched user inputs with pre-defined responses through a rule base. The limitations of these AI Agents were significant. The cost of rule maintenance was high, and a large

number of rules needed to be manually defined, making it difficult to cover complex scenarios. The generalization ability was weak, and they could not handle inputs not pre-defined in the rules, lacking flexibility.

### 3.2 Rise of Specialized Agents

With the introduction of machine learning technology, AI Agents began to decouple into three modules: perception, decision-making, and execution, and supported function expansion through a plug-in mechanism. The perception module integrates models such as speech recognition, optical character recognition, and natural language processing, converting multi-modal inputs into structured data. The decision-making module generates action sequences based on reinforcement learning or planning algorithms. The execution module calls external application programming interfaces or hardware devices to complete tasks. At this stage, AI Agents still relied on manual design processes and had limited autonomy.

### 3.3 Large Model-driven Period

The proposal of the Transformer architecture promoted the development of large language models, and AI Agents entered the large model-driven stage. Large language models, through large-scale parameterization and massive data training, possess strong language understanding and generation capabilities. For example, GPT series models enable AI Agents to conduct multi-round conversations and task decomposition. Through human feedback reinforcement learning, the quality of the conversation is optimized, promoting the evolution of AI Agents towards general intelligent agents. At this stage, AI Agents began to have initial autonomy, but still required manual design of tool invocation logic.

### 3.4 Agent Explosion Period

With the continuous optimization of the Transformer architecture and the maturity of AI Agent frameworks, AI Agents entered the explosion period. Multi-agent collaboration decomposes complex tasks through master-slave agents or peer-agent architectures. For example, one agent is responsible for planning, and another for execution. Autonomous evolution continuously optimizes the model and strategies through lifelong learning. For instance, AI Agents continuously learn new knowledge during operation and adapt to environmental changes. Edge computing deployment deploys lightweight agents to terminal devices to reduce latency and reliance on cloud services. For example, AI Agents in smart home systems process user instructions locally, improving response speed.

## 4. Future Outlook: Challenges and Opportunities of AI Agents

### 4.1 Technical Challenges

Currently, the large models relied upon by AI Agents face numerous technical challenges. One of the major issues is interpretability. The decision-making process of large models is like a "black box", lacking transparency, which limits their application in fields that have strict requirements for decision-making basis. Therefore, an interpretability framework needs to be established, such as visualizing attention weights to clearly present the reasoning path of AI Agents, allowing people to understand their decision-making process. Data bias is also an issue that cannot be ignored. The training data may contain social biases, and if not addressed, it will lead to unfair results from the model. Therefore, measures such as data cleaning and algorithm optimization should be taken to reduce the impact of bias on the fairness of the model. Additionally, the demand for computing resources is a bottleneck for the application of large models. Large model training and inference require huge computing resources, which are not only costly but also may not be deployed on a large scale due to hardware limitations. Thus, exploring model compression and quantization techniques to reduce inference latency and improve computational efficiency has become an urgent issue to be addressed.

### 4.2 Application Opportunities

AI Agents demonstrate significant application potential in multiple fields. In the field of intelligent manufacturing, they can achieve intelligent upgrades of production lines by using customized AI Agents and transfer learning techniques to precisely optimize process flows, improving production efficiency and product quality. In the field of intelligent healthcare, AI Agents can assist doctors in disease diagnosis and treatment plan formulation, combining medical images and medical record texts to provide personalized diagnosis and treatment suggestions for patients, enhancing the accuracy and effectiveness of medical services. In the field of intelligent transportation, AI Agents can predict traffic congestion through real-time data analysis, dynamically adjust signal light timing, optimize traffic flow management, improve road traffic efficiency, and alleviate urban traffic pressure.

## 5. Conclusion

From the foundation of machine learning to the breakthrough of Transformer, the evolution of AI Agents reflects the leap of artificial intelligence technology from specialized tools to general intelligent agents. Machine learning provides AI Agents with basic cognitive abilities, while the Transformer architecture endows them with powerful reasoning and decision-making capabilities. In the future, with the continuous advancement of technology and the expansion of application scenarios, AI Agents will play an important role in more fields and become digital companions of human society.

## References

- [1] Learning, D. (2020). Deep learning. High-dimensional fuzzy clustering.
- [2] Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). Deep learning (Vol. 1, No. 2, pp. 1-800). Cambridge: MIT press.
- [3] Jain, A. K., Murty, M. N., & Flynn, P. J. (1999). Data clustering: a review. *ACM computing surveys (CSUR)*, 31(3), 264-323.
- [4] Barto, A. G. (2021). Reinforcement learning: An introduction. by richard' s sutton. *SIAM Rev*, 6(2), 423.
- [5] Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. (2009). The elements of statistical learning: data mining, inference, and prediction (Vol. 2, pp. 1-758). New York: springer.
- [6] Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2017). ImageNet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6), 84-90.
- [7] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- [8] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... & Hassabis, D. (2015). Human-level control through deep reinforcement learning. *nature*, 518(7540), 529-533.
- Xu, L., Zhang, X., & Dong, Q. (2020). CLUECorpus2020: A large-scale Chinese corpus for pre-training language model. *arXiv preprint arXiv:2003.01355*.