

STEVE Series: Step-by-Step Construction of Agent Systems in Minecraft

Anonymous CVPR submission

Paper ID 2

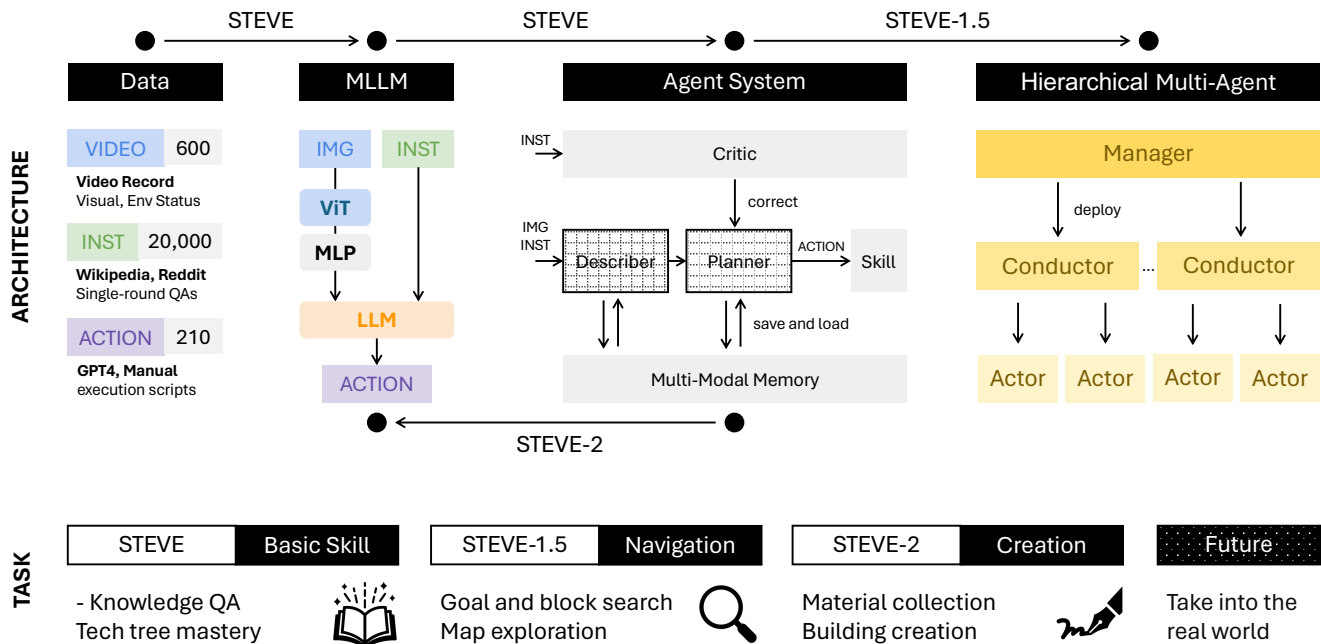


Figure 1. STEVE Series overview.

Abstract

Building an embodied agent system with a large language model (LLM) as its core is a promising direction. Due to the significant costs and uncontrollable factors associated with deploying and training such agents in the real world, we have decided to begin our exploration within the Minecraft environment. Our STEVE Series agents can complete basic tasks in a virtual environment and more challenging tasks such as navigation and even creative tasks, with an efficiency far exceeding previous state-of-the-art methods by a factor of $2.5\times$ to $7.3\times$. We begin our exploration with a vanilla large language model, augmenting it with a vision encoder and an action codebase trained on our collected high-quality dataset STEVE-21K. Subsequently, we enhanced it with a Critic and memory to transform it into a complex system. Finally, we constructed a hierarchical multi-agent system. Our recent work explored

how to prune the agent system through knowledge distillation. In the future, we will explore more potential applications of STEVE agents in the real world. The code, data, and models are available at [site](#).

1. Data and Environment

The STEVE-21K dataset is integral for training the multi-modal Large Language Models (LLMs) in the STEVE Series, containing 600 Vision-Environment pairs, 20,000 Question-Answering pairs, and 210 Skill-Code pairs to enhance agents' interaction and task execution in Minecraft. Our simulation environment utilizes MineDojo [1] and Mineflayer [5] APIs, providing a realistic setting for high-fidelity agent performance.

2. Multi-Modal LLMs

The STEVE Series advances through the integration of Multi-Modal Large Language Models (MLMs), essential

Knowledge QA		Tech Tree Mastery	
Model	preference (\uparrow)	Method	# iters (\downarrow)
Llama2-13B [8]	6.89	AutoGPT [7]	107
GPT-4 [4]	8.04	Voyager [9]	35
STEVE-13B [11]	8.12	STEVE-1 [11]	33

Table 1. **Comparison on Basic Skill.** Models preference rated 0-10 on knowledge QA and # iters stand for average iterations for task fulfillment.

for enhancing agent interactions within Minecraft. From **STEVE-1 [11]**, using the fine-tuned STEVE-13B model, to **STEVE-2 [12]** which incorporates advanced visual models like LLaVA [2, 3], each version progressively enhances the agents’ multimodal processing abilities.

3. Hierarchical Multi-Agent System

Introduced in **STEVE-1.5**, our Hierarchical Multi-Agent System enhances multi-agent cooperation for complex navigation and creation tasks in Minecraft. This system supports centralized planning and decentralized execution, enabling agents to adjust strategies and dynamically improve interaction with the environment. **STEVE-2** extends this system’s capabilities, accommodating a broader range of activities and pushing the boundaries of autonomous multi-agent systems.

4. Distill Embodied Agent into a Single Model

STEVE-2 [12] introduces a hierarchical knowledge distillation process that refines the alignment of tasks across various granularity levels within our agent system. This process incorporates the extra expert to enhance the teacher model with prior knowledge, significantly improving training quality for complex tasks. By distilling capabilities into a single model, **STEVE-2 [12]** achieves operational simplicity and superior performance, setting a new benchmark in autonomous agent capabilities within Minecraft.

5. Experiments

5.1. Basic Skill

The **STEVE series** demonstrates prowess in Knowledge Question and Answering and Tech Tree Mastery. STEVE-13B excels in producing precise Minecraft-related answers, surpassing both LLaMA2 [8] and GPT-4 [4]. In Tech Tree Mastery, **STEVE-1 [11]** progresses through Minecraft’s tech levels faster than competitors like AutoGPT [7] and Voyager [9], showcasing effective use of its vision unit to handle complex crafting tasks.

Method	# LLMs	Goal Search	Map Explore
		success (\uparrow)	# area (\uparrow)
Voyager [9]	12 / 20	64%	755
STEVE-1 [11]	20 / 24	64%	696
STEVE-2 [12]	5 / 8	91%	1493

Table 2. **Comparison on Navigation.** We list the success rate of Goal Search. # area is the average squares of blocks over 5 iterations. We list the best performance with the number of LLMs for different tasks.

Method	# LLMs	Material Collection	Building Creation
		completion (\uparrow)	FID (\downarrow)
Voyager [9]	4	72%	256.75
Creative Agents [10]	4	-	68.32
STEVE-2 [12]	8 / 2	99%	21.12

Table 3. **Comparison on Creation.** We list task completion rates and average FID scores for image quality. We list the best performance with the number of LLMs for different tasks.

5.2. Navigation

STEVE-2 [12] excels in multi-modal goal search, continuous block search, and map exploration, outperforming existing models by substantial margins. In multi-modal goal search, STEVE-2 identifies goals using various sensory inputs with performance $5.5 \times$ better than leading LLM-based methods. For map exploration, STEVE-2 updates and expands game maps with $1.9 \times$ the efficiency of previous models, using a dynamic strategy tailored to unexplored areas.

5.3. Creation

In creation tasks, **STEVE-2 [12]** significantly outperforms in material collection and building creation. It improves material gathering efficiency by $19 \times$ over Voyager [9]. Additionally, using a finetuned VQ-VAE [6] for 3D occupancy generation, STEVE-2 enhances the quality of construction, achieving a $3.2 \times$ increase in FID scores and surpassing other models and human evaluations in creative task performance.

6. Conclusion

The **STEVE series** has achieved substantial progress in multi-modal and hierarchical agent systems within Minecraft, excelling in tasks of basic skill, navigation, and creation.

Future Work The next goal is to adapt the **STEVE series’** sophisticated agent technologies for practical applications in complex, dynamic real-world environments.

095 **References**

- 096 [1] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar,
097 Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang,
098 Yuke Zhu, and Anima Anandkumar. Minedojo: Build-
099 ing open-ended embodied agents with internet-scale knowl-
100 edge. *Advances in Neural Information Processing Systems*,
101 35:18343–18362, 2022. 1
- 102 [2] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee.
103 Improved baselines with visual instruction tuning, 2023. 2
- 104 [3] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.
105 Visual instruction tuning. In *NeurIPS*, 2023. 2
- 106 [4] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv: Arxiv-*
107 *2303.08774*, 2023. 2
- 108 [5] PrismarineJS. Prismarinejs/mineflayer: Create minecraft
109 bots with a powerful, stable, and high level javascript api.,
110 2013. 1
- 111 [6] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Gener-
112 ating diverse high-fidelity images with vq-vae-2. *Advances*
113 *in neural information processing systems*, 32, 2019. 2
- 114 [7] Significant-Gravitas. Auto-gpt. [https://github.com/](https://github.com/Significant-Gravitas/Auto-GPT)
115 [Significant-Gravitas/Auto-GPT](https://github.com/Significant-Gravitas/Auto-GPT), 2023. 2
- 116 [8] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert,
117 Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov,
118 Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al.
119 Llama 2: Open foundation and fine-tuned chat models. *arXiv*
120 *preprint arXiv:2307.09288*, 2023. 2
- 121 [9] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar,
122 Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandku-
123 mar. Voyager: An open-ended embodied agent with large
124 language models. *arXiv preprint arXiv:2305.16291*, 2023. 2
- 125 [10] Chi Zhang, Penglin Cai, Yuhui Fu, Haoqi Yuan, and
126 Zongqing Lu. Creative agents: Empowering agents
127 with imagination for creative tasks. *arXiv preprint*
128 *arXiv:2312.02519*, 2023. 2
- 129 [11] Zhonghan Zhao, Wenhao Chai, Xuan Wang, Li Boyi,
130 Shengyu Hao, Shidong Cao, Tian Ye, Jenq-Neng Hwang,
131 and Gaoang Wang. See and think: Embodied agent in virtual
132 environment. *arXiv preprint arXiv:2311.15209*, 2023. 2
- 133 [12] Zhonghan Zhao, Ke Ma, Wenhao Chai, Xuan Wang, Kewei
134 Chen, Dongxu Guo, Yanting Zhang, Hongwei Wang, and
135 Gaoang Wang. Do we really need a complex agent system?
136 distill embodied agent into a single model. *arXiv preprint*
137 *arXiv:2404.04619*, 2024. 2