

AutoGLM: Autonomous Foundation Agents for GUIs

Xiao Liu¹², Bo Qin^{1†}, Dongzhu Liang^{1†}, Guang Dong^{1†}, Hanyu Lai^{12*†}, Hanchen Zhang^{12*†}, Hanlin Zhao^{1†}, Iat Long Iong^{12*†}, Jiadai Sun^{1†}, Jiaqi Wang^{1†}, Junjie Gao^{1†}, Junjun Shan^{1†}, Kangning Liu^{1†}, Shudan Zhang^{12*†}, Shuntian Yao^{1†}, Siyi Cheng^{1†}, Wentao Yao^{12*†}, Wenyi Zhao^{1†}, Xinghan Liu^{12*†}, Xinyi Liu^{1†}, Xinying Chen^{1†}, Xinyue Yang^{1†}, Yang Yang^{1†}, Yifan Xu^{12*†}, Yu Yang^{1†}, Yujia Wang^{1†}, Yulin Xu^{1†}, Zehan Qi^{12*†}, Yuxiao Dong², Jie Tang²

Project Page: <https://xiao9905.github.io/AutoGLM>

Zhipu AI¹

Tsinghua University²



Abstract

We present AUTOGLM, a new series in the ChatGLM family [11], designed to serve as foundation agents for autonomous control of digital devices through Graphical User Interfaces (GUIs). While foundation models excel at acquiring human knowledge, they often struggle with decision-making in dynamic real-world environments, limiting their progress toward artificial general intelligence. This limitation underscores the importance of developing foundation agents capable of learning through autonomous environmental interactions by reinforcing existing models. Focusing on Web Browser and Phone as representative GUI scenarios, we have developed AUTOGLM as a practical foundation agent system for real-world GUI interactions. Our approach integrates a comprehensive suite of techniques and infrastructures to create deployable agent systems suitable for user delivery. Through this development, we have derived two key insights: First, the design of an appropriate "intermediate interface" for GUI control is crucial, enabling the separation of planning and grounding behaviors, which require distinct optimization for flexibility and accuracy respectively. Second, we have developed a novel progressive training framework that enables self-evolving online curriculum reinforcement learning for AUTOGLM. Our evaluations demonstrate AUTOGLM's effectiveness across multiple domains. For web browsing, AUTOGLM achieves a 55.2% success rate on VAB-WebArena-Lite (improving to 59.1% with a second attempt) and 96.2% on OpenTable evaluation tasks. In Android device control, AUTOGLM attains a 36.2% success rate on AndroidLab (VAB-Mobile) and 89.7% on common tasks in popular Chinese APPs. Select AUTOGLM capabilities are now available through the [Qingyan Browser Plugin](#) for web applications and via [Form Applications](#) for invited Android testing. Additional results and materials will be released at <https://github.com/THUDM/AutoGLM>.

1 Introduction

Foundation models, including Large Language Models (LLMs) [5; 27; 7; 2; 42; 11] and Large Multimodal Models (LMMs) [20; 25; 26; 1], have captured widespread attention for their remarkable language understanding and generation capabilities. Through extensive self-supervised [22] pre-training on internet-scale corpora, these models have acquired not only knowledge and language

*Work done while these authors interned at Zhipu AI.

†These authors are listed alphabetically by first names.

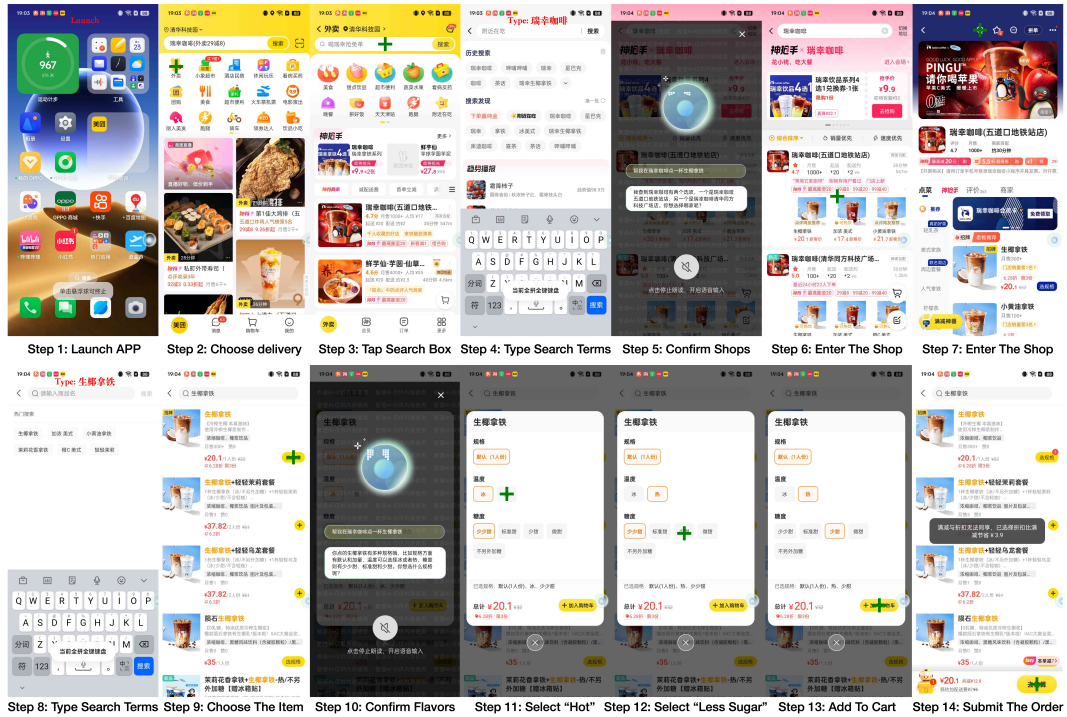


Figure 1: AUTOGLM real phone use demonstration for instruction “Order a hot coconut latte from Luckin Coffee with half sugar” for Chinese Android APP Meituan. [See more videos.](#)

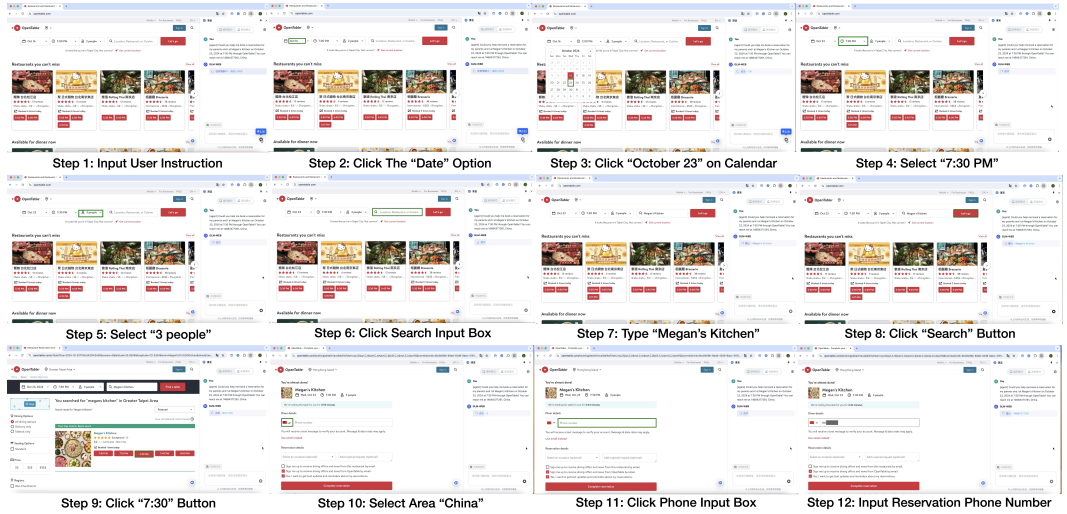


Figure 2: AUTOGLM real web-browser use demonstration for instruction “Could you help me book a reservation for my parents and I at Megan’s Kitchen on October 23, 2024 at 7:30 PM through OpenTable? You can reach me at 146xxxxxxx, China.” on the website OpenTable. [See more videos.](#)

abilities but also human-like reasoning and planning capabilities, giving rise to **LLMs as Agents** [21; 28]. These agents have demonstrated their utility across diverse domains, including coding [35; 16; 44], data analysis [14; 21], and gaming [34; 18], charting a promising course toward Artificial General Intelligence (AGI) through the development of multimodal **Foundation Agents** [23] that serve as generalists across multiple tasks and environments.

The ubiquity of digital devices presents a unique opportunity for GUI-capable agents [13; 46; 43; 17]. This domain offers several advantages: GUI simulators can be readily deployed in parallel for data

annotation and online reinforcement learning (RL); GUI environments provide rich textual and visual inputs essential for foundation model agents, but in a safer and controllable environments compared to embodied environments; and GUI agents hold broad practical appeal given their extensive potential user base. Their successful development could fundamentally transform human-device interaction.

However, the development of foundation agents for GUI faces a critical challenge: the scarcity of decision-making data in existing pre-training sets. While the internet contains vast human knowledge, it primarily consists of *static* information that inadequately captures human decision-making and environmental interaction. Building capable foundation agents requires enriching them with *dynamic* knowledge, either through direct interaction with real-world environments or through learning from synthesized trajectories. Such foundation agents can then self-evolve in the digital world, iteratively improving to achieve genuine general intelligence.

Crucially, these systems must be developed with progressive user deployment in mind. Autonomous agents are designed to augment, not replace, human capabilities. User deployment serves the dual purpose of teaching agents effective human assistance while allowing humans to adapt to intelligent assistants. This approach also enables researchers to systematically understand, discover, and examine both the potential benefits and risks of autonomous foundation agents during development.

In response to these opportunities and challenges, we introduce AUTOGLM, a series of foundation agents built upon the ChatGLM [11] model family. AUTOGLM represents a pioneering attempt to develop foundation agent prototypes for two fundamental GUI scenarios: Web Browser and Android. To address the data scarcity challenge, we employ a comprehensive suite of training techniques and develop key infrastructures for user deployment. This process has yielded two crucial insights:

- **Intermediate Interface Design:** We find it essential to design an intermediate interface that disentangles planning and grounding behaviors in foundation GUI agents. They present distinct requirements – planning demands flexibility and error recovery, while grounding emphasizes action accuracy. Their separation enables more agile development and enhanced performance.
- **Self-Evolving Online Curriculum RL [30]:** We recognize that error recovery [23] is crucial for robust and deployable agent applications, yet it remains difficult to acquire through offline training alone. Additionally, the shortage of instructions and trajectories impedes training progress. We address this challenge through self-evolving RL, implemented according to a progressive weak-to-strong curriculum schedule in an online manner.

Building on these insights, AUTOGLM demonstrates exceptional capabilities across various benchmarks and real-world tests. In Web Browsing, AUTOGLM achieves a task success rate (SR) of 55.2% on the challenging VAB-WebArena-Lite [47; 23], substantially surpassing GPT-4o’s 18.2%. With a second attempt opportunity, this improves to 59.1%. On OpenTable real-world booking tasks, AUTOGLM achieves 96.2% SR, outperforming both GPT-4o (62.6% SR) and Agent Q [29] (81.7%). Select AUTOGLM web capabilities are publicly available via the [Qingyan Browser Plugin](#) on both Chrome and Edge Plugin Store. See the real example in Figure 2.

For Android control, AUTOGLM achieves 36.2% SR on AndroidLab [37] (previously known as VAB-Mobile [23]), a comprehensive interactive Android evaluation framework. This performance exceeds both GPT-4o (31.2% SR) and Claude-3.5-Sonnet (29.0% SR). We have also implemented a practical Android application via AccessibilityService for autonomous device control. In human evaluation, AUTOGLM achieves an impressive 89.7% SR on common tasks (e.g., “Please Order a large iced Americano with half sugar from the nearest coffee shop for delivery to my company”) across popular Chinese APPs. The Android client is currently available for invited internal testing through [Form Applications](#). See the real example in Figure 1.

2 AUTOGLM: Techniques and Insights

In this section, we will give an overview of the techniques involved in developing AUTOGLM. Particularly, we will discuss the two important insights that enable AUTOGLM’s significant improvements compared to existing LLM or LMM-based GUI Agents.

2.1 Important Techniques

Training agents can be different from training ordinary LLMs or LMMs. A key obstacle lies in the lack of high-quality trajectory data that entails the decision-making process. Following are some useful techniques we realized during the project.

Pre-training. Generally, there is little agent-related data on internet text corpora, making LLMs fall short of effectively act as agents. Additionally, existing LMM pre-training, which is primarily “visual instruction tuning”, models the alignment between texts and images without sufficiently learning from sequential multimodal data [4; 10]. As a result, properly leveraging existing online data with weak-supervised decision-making signals in pre-training would actually help. Besides, for multimodal perception, high-resolution visual inputs are very important according to CogAgent [13] and our observations, especially when using grounding strategies like Set-of-Marks (SoM) prompting [38].

Large Multimodal Models (LMMs). LMMs are important for GUI understanding and manipulation. Traditionally in Robotic Process Automation (RPA), the paradigm has been using Optical Character Recognition (OCR) catcher to match key elements in human handcrafted automating programs, which cannot be scaled and generalized. LMMs, instead, can perform fuzzy matching and do long-horizon planning thanks to its strong grasping of commonsense and GUI environments from pre-training. Nevertheless, LMMs still require much training to gain strong planning and reasoning abilities necessary for agent tasks.

Behavior Cloning (Supervised Fine-tuning). Behavior Cloning (BC) is a key strategy for training agents from scratch with high-quality expert trajectories. The strategy has also been verified as effective for LLM and LMM-based agent training [24; 41; 6; 13; 17; 23]. Nevertheless, it is of extreme cost and time to collect expert trajectories. Moreover, a fundamental problem of using BC is that agents only learn to imitate experts’ behaviors step-by-step without fully understanding its goal. When expert trajectories are oracle (mostly the case for maintaining training stability), agents fail to foster abilities to recover from errors well [23].

Curriculum Learning. Agent tasks are usually of substantially varied difficulties. As a result, it is wise to progressively add difficulty to training with a curriculum schedule. For example, AutoWebGLM [17; 15] adopts a multi-stage curriculum, where agent models are sequentially trained with single-step tasks, simple few-step tasks, and complicated long-horizon tasks. DigiRL [3] also proposes a simple curriculum to filter appropriate tasks from a fixed set of instructions according to the corresponding agent capabilities at a certain timestamp. We basically find the strategy very useful for building foundation agents with complex goal-achieving abilities.

Reward Modeling (RM). To enable online RL with foundation agents, a proper RM is necessary for providing supervision. Traditionally, many RL agents are trained with limited tasks with precise rule-based reward functions. However, foundation agents based on LLMs and LMMs are targeting generalist mission accomplishment in open worlds, which contradicts task-specific reward functions’ abilities. Therefore, it is crucial to build generalizable RMs that can deal with a wide range of real-world agent tasks. Specifically, RMs can be categorized to outcome-supervised ORM and process-supervised PRM [19; 8; 40], which provide different granularities of effective supervision.

Reinforcement Learning (RL). Compared to BC, RL from a narrow sense can better learn from failures. This is especially important for foundation agent training since high-quality expert trajectories are extraordinarily hard to acquire [24]. However, the challenge in applying RL to foundation agent training lies in the inefficiency of sampling in environments. The problem can be understood from two aspects: 1) Simulators: when agents are exploring in the Web or Android environments, their efficiency is bounded by internet connection speed and maximum degree of parallelism. Environments like Android Virtual Devices are quite memory-consuming [23]. 2) Sample Diversity: LLMs and LMMs are trained to output certain function-based actions. The strict function formatting usually requires overfitting training with the model, resulting in stubborn monotonous sampling results even when they are inferenced with a high temperature [33].

Despite the challenge, we believe to scale up RL and post-training on foundation models is crucial for building strong foundation agents, as indicated by the success of OpenAI o1. It is unlikely to build general intelligence without letting it learn from interactions with real-world environments.

Table 1: Experiments on Intermediate Interface Design on VAB-WebArena-Lite [47; 23].

Observation Type	gpt-4o (text)	gpt-4o (visual)	gpt-4-vision-preview (visual)
End-to-End Agent	14.3%	18.2%	18.8%
Intermediate Interface Design	18.1% (+3.8%)	27.3% (+9.1%)	36.4% (+17.6%)

2.2 Insight 1: Intermediate Interface Design

During the development, we find intermediate interface design vital for disentangling the behaviors of planning and grounding in foundation agents. By separating them into different modules, foundation agents can be improved from both dimensions of flexibility and accuracy without interference.

The intuition is simple: we find existing LLMs and LMMs to be more capable in planning than in grounding when executing agent tasks on existing benchmarks. While the planning could still be significantly improved, a majority of current errors arise from incorrect element identification in the grounding period [23]. For example, a typical action generated in VAB-WebArena-Lite when testing with visual inputs would be:

```
# End-to-End Agent
do(action="Click", element_coordinates=[823,684])
```

where the element “4” may refer to a "Submit" button on Reddit. If we change the format to the following:

```
# Agent with Intermediate Interface Design
# Planner
do(action="Click", element_description="the 'Submit' button on the bottom right")
# Grounder
find_coordinates_by_instruction("the 'Submit' button on the bottom right")
# Expected Output: [823,684]
```

In this way, planner and grounder abilities could be separately improved. Actually, it is much easier to harvest massive grounding data from automatic construction from unsupervised environmental observations. In our experiments (Cf. Table 1), we find the strategy with the grounder we trained to be very useful for improving proprietary LLM/LMM API-based planners. Our observation is similar to another concurrent work [12] which explores a universal grounding model for GUI agents.

2.3 Insight 2: Self-Evolving Online Curriculum RL

While Intermediate Interface Design helps alleviate the issue of inaccurate grounding, planning is still a problem. Many existing agent works in literature base their frameworks on proprietary LLM/LMM APIs, whose planning abilities consequently fail to be improved by training.

As a result, we decide to explore training in-house planners via RL. It is very challenging as it lacks a sufficient amount of either user tasks or expert trajectories. We develop a self-evolving online curriculum RL framework—WebRL [30]—for training foundation agents from scratch. Take WebArena [47] environment as an example, we adopt the actor-critic RL framework for training. Briefly speaking, we identify the most difficult issues when we apply curriculum RL to the problem—task data scarcity and policy distribution drift.

Task Data Scarcity. Leveraging around 1,000 BC data provided by VisualAgentBench [23], we initialize GLM-4-9B to 22.4% SR. At this point, we have run out of either task instructions or oracle

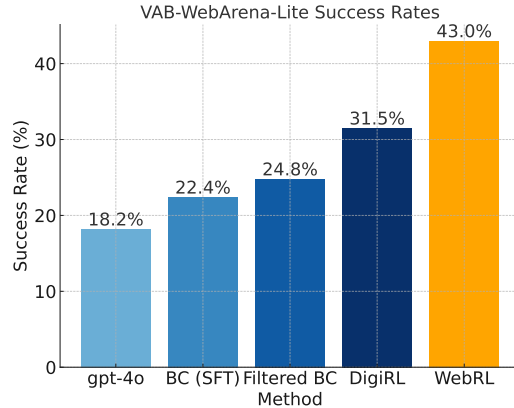


Figure 3: VAB-WebArena-Lite [47; 23], methods are all trained using GLM-4-9B-Base [11].

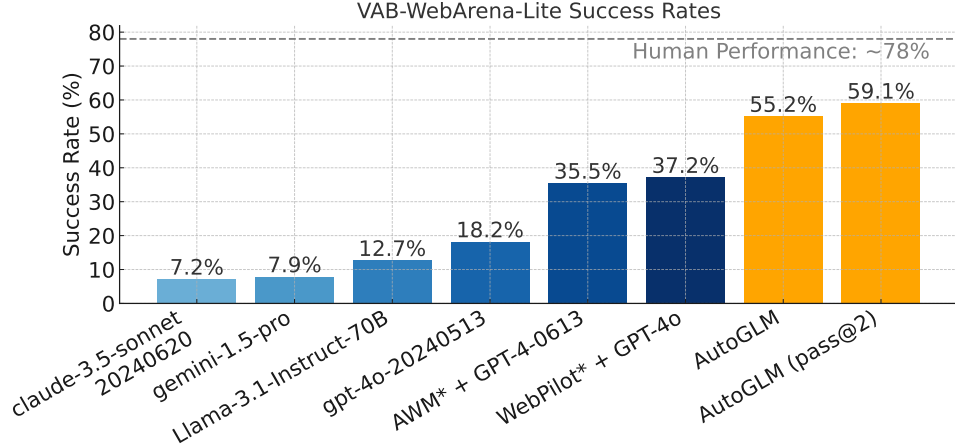


Figure 4: Success rates of different agents on VAB-WebArena-Lite. *Results of AWM and WebPilot are evaluated on full WebArena and taken from reporting.

trajectories. Thus we apply self-evolving techniques to augment failed task instructions during online roll-outs, mutating instructions to be more complicated or simpler. These self-evolved instructions are filtered by the critic and then used for roll-outs in the next iterative training phase.

Policy Distribution Drift. One significant problem of curriculum learning is the policy distribution drift during the progressive curriculum schedule. We develop a KL-constrained policy update for agent training, together with actor confidence filtered experience replay. The ablation study indicates that the designs are indispensable for consistent performance improvement during iterative training.

3 Results

In this section, we report our evaluation of AUTOGLM on both Web and Android-oriented tasks.

3.1 Evaluated on Web

We adopt three interactive benchmarks: VAB-WebArena-Lite [47; 23] and an online human evaluation dataset OpenTable [29]. AUTOGLM experiences training optimization in these environments.

VAB-WebArena-Lite [47; 23]. VAB-WebArena-Lite¹ is a refined 165-task subset of the original 812-task WebArena [47] with manual verification of answers and judge functions. Its design intention is to speed up the evaluation on WebArena and ensure judging correctness. We evaluate representative proprietary LLM/LMM APIs, open models [9], recent agent frameworks [36; 45], and AUTOGLM. Results in Figure 4 show that AUTOGLM has significantly advanced on the benchmark, narrowing the performance gap between autonomous agents and humans.

OpenTable Eval [29]. Following Agent Q [29], we also evaluate AUTOGLM on a real website OpenTable, which provides an online open booking service. Since the test set of [29] is undisclosed, we reconstruct a 200-sample test set according to the example provided in its paper (“Book reservation for restaurant Cecconi’s on OpenTable for 4 people on May 22 2024 at 7:00 PM”) and run evaluation on the real OpenTable website with human evaluation. Results are in Figure 5. AUTOGLM outperforms both gpt-4o and Agent Q on this real-world website.

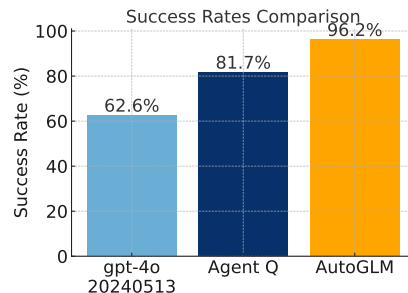


Figure 5: OpenTable Eval [29].

¹<https://github.com/THUDM/VisualAgentBench/blob/main/VAB-WebArena-Lite>

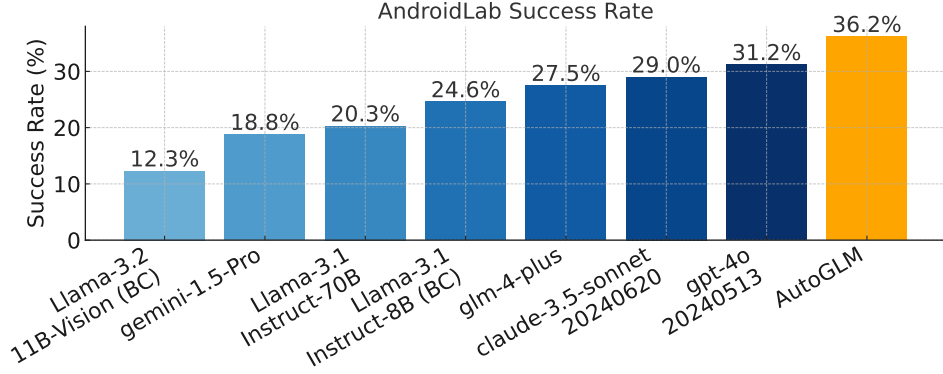


Figure 6: Success rates of different agents on AndroidLab. BC indicates behavior cloning training.

Table 2: Examples of test queries for evaluating AUTOGLM on Chinese APPs.

APP	Test Query Example
WeChat	Post a complimentary comment on Alice’s latest post on Moments
Meituan	Order a cold coconut latte from nearest coffee shop, with half sugar
Taobao	Check the shipping/delivery status of the shirt I ordered
Dianping	What’s the nearest restaurant from Beijing’s must-eat list to my location
Amap	Order a ride/taxi to Taikoo Li Sanlitun, immediately
Xiaohongshu	Bookmark and summarize the most liked travel guide post about Aranya
12306	Book a bus ticket from Zhuhai to Guangzhou for Saturday

3.2 Evaluated on Android

We evaluate AUTOGLM’s Android abilities on the academic benchmark AndroidLab [37] (i.e., VAB-Mobile [23]) and frequent tasks in common Chinese Mobile APPs on Android.

AndroidLab [37] (VAB-Mobile [23]). AndroidLab is an interactive Android benchmark and development environments that support reproducible evaluation, covering systems and some offline deployable English APPs. Compared to some existing benchmarks such as AITW [32], its interactive nature allows more practical evaluation of foundation agents for Android and improvement via RL. We evaluate representative proprietary LLM/LMM APIs, open models [9] fine-tuned on provided BC data, and AUTOGLM. Results are shown in Figure 6. AUTOGLM achieves 36.2% SR, the best-performed one among all compared agents.

Human Evaluation on Chinese Android APPs. To test the practicality of AUTOGLM being deployed for public users, we carefully examine it on frequent tasks in 7 common Chinese Android APPs, including WeChat, Meituan, Taobao, Dianping, Amap, Xiaohongshu, and 12306.

We curate a test query set (Cf. Table 2) for evaluating AUTOGLM’s real performance in the user delivery setting, where the final success rate is determined by human evaluation on the whole executing trajectories. Instead of evaluating an Android Virtual Device (AVD) as in AndroidLab and previous work [31; 39], our evaluation is conducted in physical Android phones implemented with AccessibilityService applications to reflect the practical scenarios of foundation agents for phone use. Results are shown in Figure 7. We classify results into 3 types for better understanding of AUTOGLM:

- **Success:** The task is completely successful, fulfilling all requirements in user instructions.

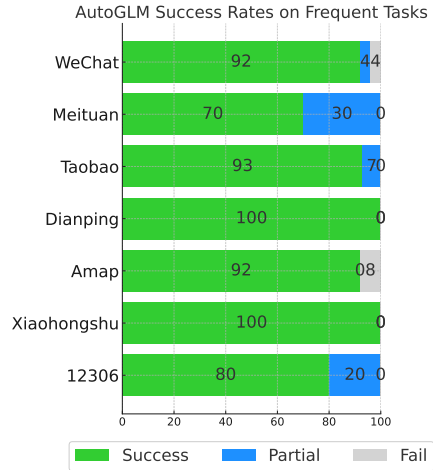


Figure 7: Human evaluated success rates of AUTOGLM on Chinese APPs.

- **Partial:** The task is partially done in the correct direction without completing some following procedures to fulfill user requirements.
- **Fail:** The task is terminated too early, gets stuck in the middle, or goes in the wrong direction.

As we observe, AUTOGLM works decently on evaluated APPs. While currently, it is unable to solve all tasks perfectly, unfinished tasks are able to be half completed, which would still be of assistance in practical scenes to users to speed up GUI operations.

4 Conclusion

Through this work, we introduced AUTOGLM, a series of foundation agents built upon the ChatGLM model family that demonstrates strong capabilities in GUI operation across web browsing and Android environments. Our key contributions include the design of an intermediate interface that effectively disentangles planning and grounding behaviors, and the development of a self-evolving online curriculum RL approach that enables robust error recovery and performance improvement. The strong empirical results across various benchmarks, including a 55.2% success rate on VAB-WebArena-Lite and 36.2% on AndroidLab, along with successful real-world deployments through browser plugins and Android applications, demonstrate AUTOGLM’s potential as a significant step toward developing practical foundation agents for GUI interaction.

References

- [1] Anthropic. Claude 3.5 sonnet, 2024.
- [2] Anthropic. Introducing the next generation of claude, 2024.
- [3] H. Bai, Y. Zhou, M. Cemri, J. Pan, A. Suhr, S. Levine, and A. Kumar. Digirl: Training in-the-wild device-control agents with autonomous reinforcement learning. *arXiv preprint arXiv:2406.11896*, 2024.
- [4] B. Baker, I. Akkaya, P. Zhokov, J. Huizinga, J. Tang, A. Ecoffet, B. Houghton, R. Sampedro, and J. Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654, 2022.
- [5] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS’20*, Red Hook, NY, USA, 2020. Curran Associates Inc.
- [6] B. Chen, C. Shu, E. Shareghi, N. Collier, K. Narasimhan, and S. Yao. Fireact: Toward language agent fine-tuning. *arXiv preprint arXiv:2310.05915*, 2023.
- [7] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [8] A. Creswell, M. Shanahan, and I. Higgins. Selection-inference: Exploiting large language models for interpretable logical reasoning. *arXiv preprint arXiv:2205.09712*, 2022.
- [9] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [10] L. Fan, G. Wang, Y. Jiang, A. Mandlekar, Y. Yang, H. Zhu, A. Tang, D.-A. Huang, Y. Zhu, and A. Anandkumar. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *Advances in Neural Information Processing Systems*, 35:18343–18362, 2022.
- [11] T. GLM, A. Zeng, B. Xu, B. Wang, C. Zhang, D. Yin, D. Rojas, G. Feng, H. Zhao, H. Lai, et al. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024.

- [12] B. Gou, R. Wang, B. Zheng, Y. Xie, C. Chang, Y. Shu, H. Sun, and Y. Su. Navigating the digital world as humans do: Universal visual grounding for gui agents. *arXiv preprint arXiv:2410.05243*, 2024.
- [13] W. Hong, W. Wang, Q. Lv, J. Xu, W. Yu, J. Ji, Y. Wang, Z. Wang, Y. Dong, M. Ding, et al. Cogagent: A visual language model for gui agents. *arXiv preprint arXiv:2312.08914*, 2023.
- [14] X. Hu, Z. Zhao, S. Wei, Z. Chai, G. Wang, X. Wang, J. Su, J. Xu, M. Zhu, Y. Cheng, et al. Infiagent-dabench: Evaluating agents on data analysis tasks. *arXiv preprint arXiv:2401.05507*, 2024.
- [15] I. L. Iong, X. Liu, Y. Chen, H. Lai, S. Yao, P. Shen, H. Yu, Y. Dong, and J. Tang. Openwebagent: An open toolkit to enable web agents on large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 72–81, 2024.
- [16] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. R. Narasimhan. Swe-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2023.
- [17] H. Lai, X. Liu, I. L. Iong, S. Yao, Y. Chen, P. Shen, H. Yu, H. Zhang, X. Zhang, Y. Dong, and J. Tang. Autowebglm: Bootstrap and reinforce a large language model-based web navigating agent, 2024.
- [18] J. Light, M. Cai, S. Shen, and Z. Hu. Avalonbench: Evaluating llms playing the game of avalon. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023.
- [19] H. Lightman, V. Kosaraju, Y. Burda, H. Edwards, B. Baker, T. Lee, J. Leike, J. Schulman, I. Sutskever, and K. Cobbe. Let’s verify step by step. *arXiv preprint arXiv:2305.20050*, 2023.
- [20] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.
- [21] X. Liu, H. Yu, H. Zhang, Y. Xu, X. Lei, H. Lai, Y. Gu, H. Ding, K. Men, K. Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*, 2023.
- [22] X. Liu, F. Zhang, Z. Hou, L. Mian, Z. Wang, J. Zhang, and J. Tang. Self-supervised learning: Generative or contrastive. *IEEE transactions on knowledge and data engineering*, 35(1):857–876, 2021.
- [23] X. Liu, T. Zhang, Y. Gu, I. L. Iong, Y. Xu, X. Song, S. Zhang, H. Lai, X. Liu, H. Zhao, et al. Visualagentbench: Towards large multimodal models as visual foundation agents. *arXiv preprint arXiv:2408.06327*, 2024.
- [24] R. Nakano, J. Hilton, S. Balaji, J. Wu, L. Ouyang, C. Kim, C. Hesse, S. Jain, V. Kosaraju, W. Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- [25] OpenAI. New models and developer products announced at devday, 2023.
- [26] OpenAI. Hello gpt-4o, 2024.
- [27] R. OpenAI. Gpt-4 technical report. *arXiv*, pages 2303–08774, 2023.
- [28] J. S. Park, J. O’Brien, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22, 2023.
- [29] P. Putta, E. Mills, N. Garg, S. Motwani, C. Finn, D. Garg, and R. Rafailov. Agent q: Advanced reasoning and learning for autonomous ai agents. *arXiv preprint arXiv:2408.07199*, 2024.
- [30] Z. Qi, X. Liu, I. L. Iong, H. Lai, X. Sun, J. Sun, X. Yang, Y. Yang, S. Yao, W. Xu, J. Tang, and Y. Dong. Webrl: Training llm web agents via self-evolving online curriculum reinforcement learning. 2024.

- [31] C. Rawles, S. Clinckemaulle, Y. Chang, J. Waltz, G. Lau, M. Fair, A. Li, W. Bishop, W. Li, F. Campbell-Ajala, et al. Androidworld: A dynamic benchmarking environment for autonomous agents. *arXiv preprint arXiv:2405.14573*, 2024.
- [32] C. Rawles, A. Li, D. Rodriguez, O. Riva, and T. Lillicrap. Android in the wild: A large-scale dataset for android device control. *Advances in Neural Information Processing Systems*, 36, 2024.
- [33] E. Wang, F. Cassano, C. Wu, Y. Bai, W. Song, V. Nath, Z. Han, S. Hendryx, S. Yue, and H. Zhang. Planning in natural language improves llm search for code generation. *arXiv preprint arXiv:2409.03733*, 2024.
- [34] G. Wang, Y. Xie, Y. Jiang, A. Mandlekar, C. Xiao, Y. Zhu, L. Fan, and A. Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.
- [35] X. Wang, B. Li, Y. Song, F. F. Xu, X. Tang, M. Zhuge, J. Pan, Y. Song, B. Li, J. Singh, et al. Opendevin: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024.
- [36] Z. Z. Wang, J. Mao, D. Fried, and G. Neubig. Agent workflow memory. *arXiv preprint arXiv:2409.07429*, 2024.
- [37] Y. Xu, X. Liu, X. Sun, S. Cheng, H. Yu, H. Lai, S. Zhang, D. Zhang, J. Tang, and Y. Dong. Androidlab: Training and systematic benchmarking of android autonomous agents. 2024.
- [38] J. Yang, H. Zhang, F. Li, X. Zou, C. Li, and J. Gao. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. *arXiv preprint arXiv:2310.11441*, 2023.
- [39] Z. Yang, J. Liu, Y. Han, X. Chen, Z. Huang, B. Fu, and G. Yu. Appagent: Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023.
- [40] E. Zelikman, Y. Wu, J. Mu, and N. Goodman. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems*, 35:15476–15488, 2022.
- [41] A. Zeng, M. Liu, R. Lu, B. Wang, X. Liu, Y. Dong, and J. Tang. Agenttuning: Enabling generalized agent abilities for llms. *arXiv preprint arXiv:2310.12823*, 2023.
- [42] A. Zeng, X. Liu, Z. Du, Z. Wang, H. Lai, M. Ding, Z. Yang, Y. Xu, W. Zheng, X. Xia, et al. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414*, 2022.
- [43] C. Zhang, Z. Yang, J. Liu, Y. Han, X. Chen, Z. Huang, B. Fu, and G. Yu. Appagent: Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023.
- [44] S. Zhang, H. Zhao, X. Liu, Q. Zheng, Z. Qi, X. Gu, X. Zhang, Y. Dong, and J. Tang. Natural-codebench: Examining coding performance mismatch on humaneval and natural user prompts. *arXiv preprint arXiv:2405.04520*, 2024.
- [45] Y. Zhang, Z. Ma, Y. Ma, Z. Han, Y. Wu, and V. Tresp. Webpilot: A versatile and autonomous multi-agent system for web task execution with strategic exploration. *arXiv preprint arXiv:2408.15978*, 2024.
- [46] B. Zheng, B. Gou, J. Kil, H. Sun, and Y. Su. Gpt-4v (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614*, 2024.
- [47] S. Zhou, F. F. Xu, H. Zhu, X. Zhou, R. Lo, A. Sridhar, X. Cheng, T. Ou, Y. Bisk, D. Fried, et al. Webarena: A realistic web environment for building autonomous agents. In *The Twelfth International Conference on Learning Representations*, 2023.