
Automated Design of Agentic Systems

Shengran Hu^{1,2}
srhu@cs.ubc.ca

Cong Lu^{1,2}
conglu@cs.ubc.ca

Jeff Clune^{1,2,3}
jclune@gmail.com

¹University of British Columbia

²Vector Institute

³Canada CIFAR AI Chair

Abstract

Researchers are investing substantial effort in developing powerful general-purpose agents, wherein Foundation Models are used as modules within *agentic systems* (e.g. Chain-of-Thought, Self-Reflection, Toolformer). However, the history of machine learning teaches us that hand-designed solutions are eventually replaced by learned solutions. We describe a newly forming research area, **Automated Design of Agentic Systems (ADAS)**, which aims to *automatically create powerful agentic system designs, including inventing novel building blocks and/or combining them in new ways*. We further demonstrate that there is an unexplored yet promising approach within ADAS where agents can be defined in code and new agents can be automatically discovered by a meta agent programming ever better ones in code. Given that programming languages are Turing Complete, this approach theoretically enables the learning of *any possible* agentic system: including novel prompts, tool use, workflows, and combinations thereof. We present a simple yet effective algorithm named Meta Agent Search to demonstrate this idea, where a meta agent iteratively programs interesting new agents based on an ever-growing archive of previous discoveries. Through extensive experiments across multiple domains including coding, science, and math, we show that our algorithm can progressively invent agents with novel designs that greatly outperform state-of-the-art hand-designed agents. Importantly, we consistently observe the surprising result that agents invented by Meta Agent Search maintain superior performance even when transferred across domains and models, demonstrating their robustness and generality. Provided we develop it safely, our work illustrates the potential of an exciting new research direction toward automatically designing ever-more powerful agentic systems to benefit humanity.

1 Introduction

Foundation Models (FMs) such as GPT [56, 54] and Claude [2] are increasingly being adopted as powerful general-purpose agents for tasks requiring flexible reasoning and planning [79]. However, solving complex real-world problems often requires agents to function as compound systems with multiple components, such as search engines, code execution, and database queries, rather than relying on a single model query [91, 65]. To address these challenges, many building blocks have been proposed, including chain-of-thought planning [83, 87, 30], memory structures [94, 39], tool use [68, 61], and self-reflection [47, 72]. While these approaches have shown success [79], they often require significant manual tuning and domain-specific expertise. History shows that manually created systems are often replaced by more efficient, learned solutions as compute and data scale [14]. Examples include the replacement of hand-crafted vision features like HOG [17] with learned CNN features [35], and the rise of AutoML and AI-Generating Algorithms (AI-GAs)[33, 14] over hand-designed AI systems. For instance, Neural Architecture Search[22, 70] now outperforms manually designed CNNs, learned loss functions [43] surpass hand-designed ones like DPO [62],

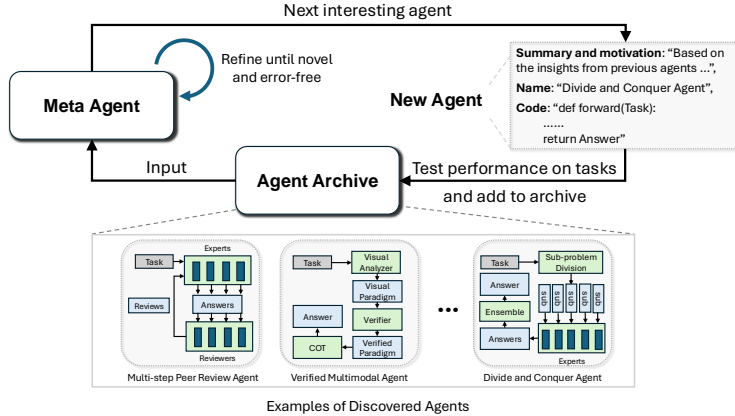


Figure 1: **Overview of the proposed algorithm Meta Agent Search and examples of discovered agents.** In our algorithm, we instruct the “meta” agent to iteratively program new agents, test their performance on tasks, add them to an archive of discovered agents, and use this archive to inform the meta agent in subsequent iterations. We show three example agents across our runs, with all names generated by the meta agent.

and AI Scientist [44] automates novel ML algorithm development. Additionally, works like OMNI-EPIC [23] automatically generate countless robotics learning environments, showing greater creativity and efficiency than manual methods.

Therefore, in this paper, we introduce a new research area called **Automated Design of Agentic Systems (ADAS)**, which aims to invent novel building blocks and design powerful agentic systems automatically (Appendix A). We argue that ADAS could be the fastest path to developing powerful agents, with initial evidence showing that learned agents can greatly outperform hand-designed ones. Given the vast number of building blocks yet to be discovered in agentic systems (Appendix B), it would take considerable time for the research community to identify them all. Even if discovered, combining these components into effective systems for real-world applications would remain challenging due to the complexity of their interactions. In contrast, ADAS allows building blocks and agents to be learned in an automated manner, potentially saving human effort and accelerating the discovery of more effective solutions. While a few existing works focus primarily on prompt design [86, 24], limiting their flexibility, we propose a novel approach to ADAS where the entire agentic system is defined in code, enabling a “meta” agent to program new agents. Since programming languages like Python are Turing Complete [5, 36], searching within a code space allows ADAS algorithms to discover *any* possible agentic systems, including components such as prompts, tool use, and workflows. With FMs becoming increasingly proficient at coding, we can use them as meta agents to create agents in code. In this paper, we present Meta Agent Search, one of the first algorithms in ADAS that enables complete design in code space (Figure 1). The core idea of Meta Agent Search is to have a meta agent iteratively create, evaluate, and archive new agents, using this archive to inform and improve future iterations, exploring novel and interesting designs similar to existing open-endedness algorithms [92, 45].

To validate the proposed approach, we evaluate Meta Agent Search on (1) the challenging ARC logic puzzle task [13], which tests the general intelligence of AI systems, (2) four popular benchmarks on reading comprehension, math, science questions, and multi-task problem solving, and (3) the transferability of discovered agents to held-out domains and models (Section 3). Our experiments show that the discovered agents significantly outperform state-of-the-art hand-designed baselines, improving F1 scores on reading comprehension by **13.6/100** and accuracy rates on math tasks by **14.4%**. Additionally, they outperform baselines by **25.9%** on GSM8K and **13.2%** on GSM-Hard after transferring across domains. These results demonstrate the potential of ADAS in automating agent design, as the discovered agents not only perform well in similar domains but also show strong performance when transferred to dissimilar domains, such as from mathematics to reading comprehension. This highlights the robustness and transferability of the agents discovered by Meta Agent Search, paving the way for further research (Section 4).

2 Our Algorithm: Meta Agent Search

We present Meta Agent Search, a simple yet effective algorithm for defining and searching for agents in code. The core idea is to use FMs as meta agents that iteratively program new agents based on an

ever-growing archive of previous discoveries. Although the meta agent can theoretically program agents from scratch, it’s inefficient to avoid providing basic functions such as FM query APIs or existing tools. Thus, we provide the meta agent with a simple framework (under 100 lines of code) that includes essential functions like querying FMs or formatting prompts. The meta agent only needs to program a ”forward” function to define new agentic systems, similar to FunSearch [67], taking task information as input and generating agent responses. As shown in Figure 1, Meta Agent Search encourages the meta agent to explore novel and interesting agents [92, 45], leveraging self-reflection [72, 47] to refine the novelty and correctness of each proposal, and performing up to three refinements when errors occur during execution. Once an agent is generated, it’s evaluated on validation data using performance metrics such as success rate or F1 score, after which it’s added to the archive. The process continues iteratively until the maximum number of iterations is reached. Detailed framework codes, agent examples, and prompts are provided in Appendices E and J.

3 Experiments

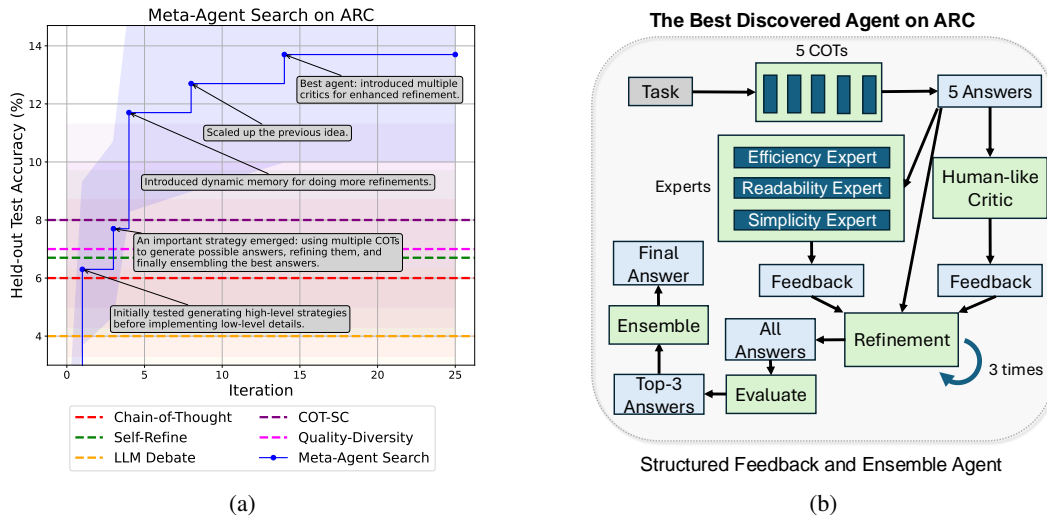


Figure 2: **The results of Meta Agent Search on the ARC challenge.** (a) Meta Agent Search progressively discovers high-performance agents based on an ever-growing archive of previous discoveries. (b) The visualization of the best agent discovered by Meta Agent Search on the ARC challenge.

We first demonstrate how Meta Agent Search discovers novel agentic systems and outperforms existing state-of-the-art hand-designed agents in the Abstraction and Reasoning Corpus (ARC) challenge [13]. Full details are listed in Appendices G and I. As shown in Figure 2a, Meta Agent Search progressively discovers agents that surpass the performance of hand-designed baselines, with key breakthroughs highlighted in the text boxes. Similar to prior work on open-endedness and AI-GAs [92, 23, 80, 81, 38], Meta Agent Search builds on a growing archive of stepping stones. For example, a critical design pattern emerged in iteration 3, where multiple COTs were used to generate, refine, and ensemble answers, becoming a foundation for later designs. The best-discovered agent (shown in Figure 2b) utilizes a complex feedback mechanism for refining answers. This mechanism evolved from ideas introduced in iterations 5, 11, and 12, focusing on feedback diversity, expert evaluations of traits like efficiency and simplicity, and simulating human-like feedback. Although these stepping stones didn’t achieve high performance immediately, later innovations successfully combined them, akin to crossover in evolution [49]. These results highlight Meta Agent Search’s ability to progressively discover agents that surpass hand-designed baselines and invent novel design patterns by combining different innovations.

Next, we investigate the potential of our algorithm to improve the capabilities of agents across math, reading, and reasoning domains, and transfer across models and domains. Experiment details are listed in Appendices D, H and I. The search results across multiple domains demonstrate that Meta Agent Search consistently discovers agents that outperform state-of-the-art hand-designed agents (Table 1). Notably, in the Reading Comprehension and Math domains, learned agents show a substantial performance gap, improving F1 scores by **13.6/100** and accuracy rates by **14.4%**, respec-

Table 1: **Performance comparison between Meta Agent Search and state-of-the-art hand-designed agents across multiple domains.** Meta Agent Search discovers superior agents compared to the baselines in every domain. The search is conducted independently for each domain.

Agent Name	F1 Score		Accuracy (%)	
	Reading Comprehension	Math	Multi-task	Science
State-of-the-art Hand-designed Agents				
Chain-of-Thought [83]	64.2 ± 0.9	28.0 ± 3.1	65.4 ± 3.3	29.2 ± 3.1
COT-SC [82]	64.4 ± 0.8	28.2 ± 3.1	65.9 ± 3.2	30.5 ± 3.2
Self-Refine [47]	59.2 ± 0.9	27.5 ± 3.1	63.5 ± 3.4	31.6 ± 3.2
LLM Debate [19]	60.6 ± 0.9	39.0 ± 3.4	65.6 ± 3.3	31.4 ± 3.2
Step-back Abstraction [95]	60.4 ± 1.0	31.1 ± 3.2	65.1 ± 3.3	26.9 ± 3.0
Quality-Diversity [45]	61.8 ± 0.9	23.8 ± 3.0	65.1 ± 3.3	30.2 ± 3.1
Role Assignment [85]	65.8 ± 0.9	30.1 ± 3.2	64.5 ± 3.3	31.1 ± 3.1
Automated Design of Agentic Systems on Different Domains				
Best Agents from Meta Agent Search	79.4 ± 0.8	53.4 ± 3.5	69.6 ± 3.2	34.6 ± 3.2

Table 2: **Performance on held-out math and non-math domains when transferring top agents from MGSM (Math).** GSM8K and GSM-Hard are the held-out math domains, while MMLU is for Multi-task, and DROP is for Reading Comprehension. Agents discovered by Meta Agent Search consistently outperform the baselines across all domains.

Agent Name	Accuracy (%)				F1 Score
	MGSM	GSM8K	GSM-Hard	MMLU	DROP
Manually Designed Agents					
Chain-of-Thought [83]	28.0 ± 3.1	34.9 ± 3.2	15.0 ± 2.5	65.4 ± 3.3	64.2 ± 0.9
COT-SC [82]	28.2 ± 3.1	37.8 ± 3.4	15.5 ± 2.5	65.9 ± 3.2	64.4 ± 0.8
Self-Refine [47]	27.5 ± 3.1	38.9 ± 3.4	15.1 ± 2.4	63.5 ± 3.4	59.2 ± 0.9
LLM Debate [19]	39.0 ± 3.4	43.6 ± 3.4	17.4 ± 2.6	65.6 ± 3.3	60.6 ± 0.9
Step-back Abstraction [95]	31.1 ± 3.2	31.5 ± 3.3	12.2 ± 2.3	65.1 ± 3.3	60.4 ± 1.0
Quality-Diversity [45]	23.8 ± 3.0	28.0 ± 3.1	14.1 ± 2.4	65.1 ± 3.1	61.8 ± 0.9
Role Assignment [85]	30.1 ± 3.2	37.0 ± 3.4	18.0 ± 2.7	64.5 ± 3.3	65.8 ± 0.9
Top Agents Searched on MGSM (Math)					
	Transferred within Math Domains		Transferred beyond Math Domains		
Dynamic Role-Playing Architecture	53.4 ± 3.5	69.5 ± 3.2	31.2 ± 3.2	62.4 ± 3.4	70.4 ± 0.9
Structured Multimodal Feedback Loop	50.2 ± 3.5	64.5 ± 3.4	30.1 ± 3.2	67.0 ± 3.2	70.4 ± 0.9
Interactive Multimodal Feedback Loop	47.4 ± 3.5	64.9 ± 3.3	27.6 ± 3.2	64.8 ± 3.3	71.9 ± 0.8

tively. While the performance gap in the Multi-task and Science domains is smaller, we hypothesize this is due to the limited knowledge in current FMs, which constrains improvement in these more challenging domains—a limitation that should diminish as FMs advance. In contrast, in Reading Comprehension and Math, where FMs have sufficient knowledge, errors stem primarily from hallucinations or calculation mistakes, which can be mitigated through well-designed agentic systems like those discovered by Meta Agent Search. Additionally, Meta Agent Search demonstrates strong transferability and generalizability. As shown in Table 2, agents discovered in math not only outperform baselines by **25.9%** on GSM8K and **13.2%** on GSM-Hard but also transfer effectively to non-math domains, outperforming state-of-the-art hand-designed agents even when not specifically designed for those domains. This generalizability also extends across different FMs (Appendix D). These results highlight Meta Agent Search’s ability to discover flexible and robust design patterns in agentic systems.

4 Safety Consideration and Conclusion

While it is highly unlikely that model-generated code will perform overtly malicious actions in our current settings and with the Foundation Models (FMs) we use, such code may still act destructively due to limitations in model capability or alignment [66, 10]. More broadly, research on more powerful AI systems raises the question of whether we should be conducting research to advance AI capabilities at all. That topic clearly includes the proposed Automated Design of Agentic Systems (ADAS) as a new area in AI-GA research, which could potentially contribute to an even faster way to create Artificial General Intelligence (AGI) than the current manual approach [14]. The question of whether and why we should pursue AGI and AI-GA has been discussed in many papers [14, 21, 4, 90, 3], and is beyond the scope of this paper. Specifically as regards ADAS, we believe it is net beneficial to publish this work. First, this work demonstrates that with the available API access to powerful FMs, it is easy to program powerful ADAS algorithms, and do so without any expensive hardware like GPUs. We feel it is beneficial to let the community know such algorithms are powerful and easy to create, so they can be informed and account for them. Moreover, by sharing this information, we hope to motivate follow-up work into safe-ADAS, such as algorithms that conduct ADAS safely during both search itself (e.g. not risking running any harmful code) and that refuse to create dishonest, unhelpful, and/or harmful agents. Such an open-source research approach to create safe-ADAS could be a better way to create safer AI systems [6, 48].

We present related work and future work in Appendices B and C. In this paper, we describe a newly forming research area, Automated Design of Agentic Systems (ADAS), which aims to *automatically invent novel building blocks and design powerful agentic systems*. We demonstrated that a promising approach to ADAS is to define agents in code, allowing new agents to be discovered by a “meta” agent programming them in code. Following this idea, we propose Meta Agent Search, where the meta agent iteratively builds on previous discoveries to program interesting new agents. The experiments show that Meta Agent Search consistently outperforms state-of-the-art hand-designed agents across an extensive number of domains, and the discovered agents transfer well across models and domains. Overall, our work illustrates the potential of an exciting new research direction toward full automation in developing powerful agentic systems from the bottom up.

Acknowledgments and Disclosure of Funding

This work was supported by the Vector Institute, the Canada CIFAR AI Chairs program, grants from Schmidt Futures and Open Philanthropy, an NSERC Discovery Grant, and a generous donation from Rafael Cosman. We thank Jenny Zhang, Rach Pradhan, Ruiyu Gou, Nicholas Ioannidis, and Eunjeong Hwang for insightful discussions and feedback.

References

- [1] Anthropic. Introducing the next generation of claude. <https://www.anthropic.com/news/claude-3-family>, March 2024. Blog post.
- [2] Anthropic. Introducing claude 3.5 sonnet. <https://www.anthropic.com/news/claude-3-5-sonnet>, June 2024. Blog post.
- [3] Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Trevor Darrell, Yuval Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, et al. Managing extreme ai risks amid rapid progress. *Science*, 384(6698):842–845, 2024.
- [4] N Bostrom. Existential Risks: analyzing human extinction scenarios and related hazards. *Journal of Evolution and Technology*, 9, 2002.
- [5] Robert S Boyer and J Strother Moore. *A mechanical proof of the Turing completeness of pure LISP*. Citeseer, 1983.
- [6] Tracey Caldwell. Ethical hackers: putting on the white hat. *Network Security*, 2011(7):10–13, 2011.
- [7] Harrison Chase. What is an agent? <https://blog.langchain.dev/what-is-an-agent/>, June 2024. Blog post.
- [8] Banghao Chen, Zhaofeng Zhang, Nicolas Langrené, and Shengxin Zhu. Unleashing the potential of prompt engineering in large language models: a comprehensive review. *arXiv preprint arXiv:2310.14735*, 2023.
- [9] Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Sesay Jaward, Karlsson Börje, Jie Fu, and Yemin Shi. Autoagents: The automatic agents generation framework. *arXiv preprint*, 2023.
- [10] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [11] Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, et al. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. In *The Twelfth International Conference on Learning Representations*, 2023.
- [12] Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E. Gonzalez, and Ion Stoica. Chatbot arena: An open platform for evaluating llms by human preference, 2024.
- [13] François Chollet. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.
- [14] Jeff Clune. Ai-gas: Ai-generating algorithms, an alternate paradigm for producing general artificial intelligence. *arXiv preprint arXiv:1905.10985*, 2019.
- [15] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [16] Antoine Cully and Yiannis Demiris. Quality and diversity optimization: A unifying modular framework. *IEEE Transactions on Evolutionary Computation*, 22(2):245–259, 2017.
- [17] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 886–893 vol. 1, 2005. doi: 10.1109/CVPR.2005.177.
- [18] Kalyanmoy Deb, Amrit Pratap, Sameer Agarwal, and TAMT Meyarivan. A fast and elitist multiobjective genetic algorithm: Nsga-ii. *IEEE transactions on evolutionary computation*, 6(2):182–197, 2002.

- [19] Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*, 2023.
- [20] Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1246.
- [21] Adrien Ecoffet, Jeff Clune, and Joel Lehman. Open questions in creating safe open-ended AI: Tensions between control and creativity. In *Conference on Artificial Life*, pages 27–35. MIT Press, 2020.
- [22] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey. *Journal of Machine Learning Research*, 20(55):1–21, 2019.
- [23] Maxence Faldor, Jenny Zhang, Antoine Cully, and Jeff Clune. Omni-epic: Open-endedness via models of human notions of interestingness with environments programmed in code. *arXiv preprint arXiv:2405.15568*, 2024.
- [24] Chrisantha Fernando, Dylan Sunil Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel. Promptbreeder: Self-referential self-improvement via prompt evolution, 2024.
- [25] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017.
- [26] Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine Learning*, pages 10764–10799. PMLR, 2023.
- [27] Ryan Greenblatt. Getting 50% sota on arc-agi with gpt-4. <https://redwoodresearch.substack.com/p/getting-50-sota-on-arc-agi-with-gpt>, July 2024. Technical Report.
- [28] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021.
- [29] Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 2023.
- [30] Shengran Hu and Jeff Clune. Thought Cloning: Learning to think while acting by imitating human thinking. *Advances in Neural Information Processing Systems*, 36, 2024.
- [31] Shengran Hu, Ran Cheng, Cheng He, Zhichao Lu, Jing Wang, and Miao Zhang. Accelerating multi-objective neural architecture search by random-weight evaluation. *Complex & Intelligent Systems*, pages 1–10, 2021.
- [32] Shihua Huang, Zhichao Lu, Kalyanmoy Deb, and Vishnu Naresh Boddeti. Revisiting residual networks for adversarial robustness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8202–8211, 2023.
- [33] Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren. *Automated machine learning: methods, systems, challenges*. Springer Nature, 2019.
- [34] Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, Heather Miller, et al. Dspy: Compiling declarative language model calls into state-of-the-art pipelines. In *The Twelfth International Conference on Learning Representations*, 2024.

- [35] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [36] Abraham Ladha. Lecture 11: Turing-completeness. <https://faculty.cc.gatech.edu/~ladha/S24/4510/L11.pdf>, 2024. CS 4510 Automata and Complexity, February 21st, 2024, Scribed by Rishabh Singhal.
- [37] LangChainAI. Langchain: Build context-aware reasoning applications. <https://github.com/langchain-ai/langchain>, 2022.
- [38] Joel Lehman and Kenneth O Stanley. Abandoning objectives: Evolution through the search for novelty alone. *Evolutionary computation*, 19(2):189–223, 2011.
- [39] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.
- [40] Fei Liu, Tong Xialiang, Mingxuan Yuan, Xi Lin, Fu Luo, Zhenkun Wang, Zhichao Lu, and Qingfu Zhang. Evolution of heuristics: Towards efficient automatic algorithm design using large language model. In *Forty-first International Conference on Machine Learning*, 2024.
- [41] Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. Dynamic llm-agent network: An llm-agent collaboration framework with agent team optimization. *arXiv preprint arXiv:2310.02170*, 2023.
- [42] Chris Lu, Sebastian Towers, and Jakob Foerster. Arbitrary order meta-learning with simple population-based evolution. In *ALIFE 2023: Ghost in the Machine: Proceedings of the 2023 Artificial Life Conference*. MIT Press, 2023.
- [43] Chris Lu, Samuel Holt, Claudio Fanconi, Alex J Chan, Jakob Foerster, Mihaela van der Schaar, and Robert Tjarko Lange. Discovering preference optimization algorithms with and for large language models. *arXiv preprint arXiv:2406.08414*, 2024.
- [44] Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The AI Scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
- [45] Cong Lu, Shengran Hu, and Jeff Clune. Intelligent go-explore: Standing on the shoulders of giant foundation models. *arXiv preprint arXiv:2405.15143*, 2024.
- [46] Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayaraman, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via coding large language models. In *The Twelfth International Conference on Learning Representations*, 2023.
- [47] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- [48] Meta. Open source ai is the path forward. <https://about.fb.com/news/2024/07/open-source-ai-is-the-path-forward/>, July 2024. News article.
- [49] Elliot Meyerson, Mark J Nelson, Herbie Bradley, Adam Gaier, Arash Moradi, Amy K Hoover, and Joel Lehman. Language model crossover: Variation through few-shot prompting. *arXiv preprint arXiv:2302.12170*, 2023.
- [50] Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su. A diverse corpus for evaluating and developing english math word problem solvers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 975–984, 2020.
- [51] Jean-Baptiste Mouret and Jeff Clune. Illuminating search spaces by mapping elites. *arXiv preprint arXiv:1504.04909*, 2015.

- [52] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- [53] Andrew Ng. Issue 253. <https://www.deeplearning.ai/the-batch/issue-253/>, June 2024. Newsletter issue.
- [54] OpenAI. Introducing chatgpt. <https://openai.com/index/chatgpt/>, November 2022. Blog post.
- [55] OpenAI. Simple evals, 2023. URL <https://github.com/openai/simple-evals>. Accessed: 2024-08-10.
- [56] OpenAI. Gpt-4 technical report, 2024.
- [57] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael 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.
- [58] Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.168.
- [59] Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. Communicative agents for software development. *arXiv preprint arXiv:2307.07924*, 2023.
- [60] Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, Zhiyuan Liu, and Maosong Sun. Scaling large-language-model-based multi-agent collaboration. *arXiv preprint arXiv:2406.07155*, 2024.
- [61] Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. Tool learning with large language models: A survey. *arXiv preprint arXiv:2405.17935*, 2024.
- [62] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [63] David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark, 2023.
- [64] Toran Bruce Richards. Autogpt. <https://github.com/Significant-Gravitas/AutoGPT>, 2023. GitHub repository.
- [65] Tim Rocktäschel. *Artificial Intelligence: 10 Things You Should Know*. Seven Dials, September 2024. ISBN 978-1399626521.
- [66] Md Omar Faruk Rokon, Risul Islam, Ahmad Darki, Evangelos E Papalexakis, and Michalis Faloutsos. SourceFinder: Finding malware Source-Code from publicly available repositories in GitHub. In *23rd International Symposium on Research in Attacks, Intrusions and Defenses (RAID 2020)*, pages 149–163, 2020.
- [67] Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.

- [68] Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=Yacmpz84TH>.
- [69] Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, et al. The prompt report: A systematic survey of prompting techniques. *arXiv preprint arXiv:2406.06608*, 2024.
- [70] Xuan Shen, Yaohua Wang, Ming Lin, Yilun Huang, Hao Tang, Xiuyu Sun, and Yanzhi Wang. Deepmad: Mathematical architecture design for deep convolutional neural network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6163–6173, 2023.
- [71] Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. Language models are multilingual chain-of-thought reasoners. In *The Eleventh International Conference on Learning Representations*, 2023.
- [72] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2023.
- [73] Kenneth O Stanley and Joel Lehman. *Why greatness cannot be planned: The myth of the objective*. Springer, 2015.
- [74] Kenneth O Stanley, Jeff Clune, Joel Lehman, and Risto Miikkulainen. Designing neural networks through neuroevolution. *Nature Machine Intelligence*, 1(1):24–35, 2019.
- [75] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [76] Sai Vemprala, Rogerio Bonatti, Arthur Bucker, and Ashish Kapoor. Chatgpt for robotics: Design principles and model abilities. Technical Report MSR-TR-2023-8, Microsoft, February 2023. URL <https://www.microsoft.com/en-us/research/publication/chatgpt-for-robotics-design-principles-and-model-abilities/>.
- [77] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv: Arxiv-2305.16291*, 2023.
- [78] Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dhharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*, 2016.
- [79] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345, 2024.
- [80] Rui Wang, Joel Lehman, Jeff Clune, and Kenneth O. Stanley. Poet: open-ended coevolution of environments and their optimized solutions. In *Proceedings of the Genetic and Evolutionary Computation Conference, GECCO '19*, page 142–151, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450361118. doi: 10.1145/3321707.3321799.
- [81] Rui Wang, Joel Lehman, Aditya Rawal, Jiale Zhi, Yulun Li, Jeffrey Clune, and Kenneth Stanley. Enhanced poet: Open-ended reinforcement learning through unbounded invention of learning challenges and their solutions. In *International conference on machine learning*, pages 9940–9951. PMLR, 2020.
- [82] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2023.

- [83] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [84] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023.
- [85] Benfeng Xu, An Yang, Junyang Lin, Quan Wang, Chang Zhou, Yongdong Zhang, and Zhen-dong Mao. Expertprompting: Instructing large language models to be distinguished experts. *arXiv preprint arXiv:2305.14688*, 2023.
- [86] Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. Large language models as optimizers. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Bb4VGOWELI>.
- [87] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=WE_vluYUL-X.
- [88] Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montserrat Gonzalez Arenas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplik, et al. Language to rewards for robotic skill synthesis. In *Conference on Robot Learning*, pages 374–404. PMLR, 2023.
- [89] Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan, Dongsheng Li, and Deqing Yang. Evo-agent: Towards automatic multi-agent generation via evolutionary algorithms. *arXiv preprint arXiv:2406.14228*, 2024.
- [90] Eliezer Yudkowsky et al. Artificial Intelligence as a positive and negative factor in global risk. *Global catastrophic risks*, 1(303):184, 2008.
- [91] Matei Zaharia, Omar Khattab, Lingjiao Chen, Jared Quincy Davis, Heather Miller, Chris Potts, James Zou, Michael Carbin, Jonathan Frankle, Naveen Rao, and Ali Ghodsi. The shift from models to compound ai systems. <https://bair.berkeley.edu/blog/2024/02/18/compound-ai-systems/>, 2024.
- [92] Jenny Zhang, Joel Lehman, Kenneth Stanley, and Jeff Clune. OMNI: Open-endedness via models of human notions of interestingness. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=AgM3MzT99c>.
- [93] Shaokun Zhang, Jieyu Zhang, Jiale Liu, Linxin Song, Chi Wang, Ranjay Krishna, and Qingyun Wu. Offline training of language model agents with functions as learnable weights. In *Forty-first International Conference on Machine Learning*, 2024.
- [94] Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. A survey on the memory mechanism of large language model based agents. *arXiv preprint arXiv:2404.13501*, 2024.
- [95] Huaixiu Steven Zheng, Swaroop Mishra, Xinyun Chen, Heng-Tze Cheng, Ed H Chi, Quoc V Le, and Denny Zhou. Take a step back: Evoking reasoning via abstraction in large language models. *arXiv preprint arXiv:2310.06117*, 2023.
- [96] Pei Zhou, Jay Pujara, Xiang Ren, Xinyun Chen, Heng-Tze Cheng, Quoc V Le, Ed H Chi, Denny Zhou, Swaroop Mishra, and Huaixiu Steven Zheng. Self-discover: Large language models self-compose reasoning structures. *arXiv preprint arXiv:2402.03620*, 2024.
- [97] Wangchunshu Zhou, Yixin Ou, Shengwei Ding, Long Li, Jialong Wu, Tiannan Wang, Jiamin Chen, Shuai Wang, Xiaohua Xu, Ningyu Zhang, et al. Symbolic learning enables self-evolving agents. *arXiv preprint arXiv:2406.18532*, 2024.
- [98] Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen Schmidhuber. Gptswarm: Language agents as optimizable graphs. In *Forty-first International Conference on Machine Learning*, 2024.

Supplementary Material

Table of Contents

A Automated Design of Agentic Systems (ADAS)	13
B Related Work	13
C Future Work	14
D Generalizability and Transferability	15
E Prompts	16
F Framework Code	19
G Experiment Details for ARC Challenge	22
H Experiment Details for Reasoning and Problem-Solving Domains	25
I Baselines	27
J Example Agents	27
K Cost of Experiments	30

A Automated Design of Agentic Systems (ADAS)

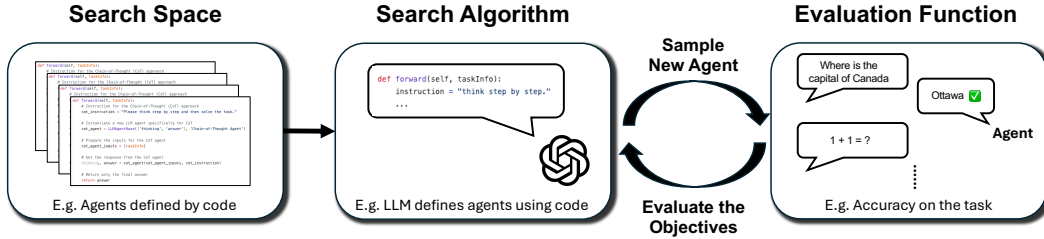


Figure 3: **The three key components of Automated Design of Agentic Systems (ADAS).** The search space determines which agentic systems can be represented in ADAS. The search algorithm specifies how the ADAS method explores the search space. The evaluation function defines how to evaluate a candidate agent on target objectives such as performance.

At the time of writing, the community has not reached a consensus on the definitions or terminologies of agents. Here, by agents we refer to agentic systems that involve Foundation Models (FMs) as modules in the workflow to solve tasks by planning, using tools, and carrying out multiple, iterative steps of processing [7, 53]. In this paper, we describe a newly forming research area Automated Design of Agentic Systems (ADAS). Similar to research areas in AI-GAs [14] and AutoML [33], such as Neural Architecture Search [22], we formulate ADAS as an optimization process and identify three key components of ADAS algorithms (Figure 3):

The **search space** in ADAS defines which agentic systems can be represented and discovered. For example, PromptBreeder [24] mutates only text prompts, leaving other components, such as workflows, unchanged, limiting the discovery of agents with different workflows. Other works explore search spaces like graph structures [98] and feed-forward networks [41]. The **search algorithm** dictates how ADAS explores the search space, balancing the exploration-exploitation trade-off [75] to quickly find high-performance systems without getting stuck in local optima. Approaches include using Reinforcement Learning [98] or iteratively generating solutions via Foundation Models (FMs) [24]. Finally, the **evaluation function** assesses candidate agents on objectives like performance, cost, latency, or safety, with accuracy on validation data being a common metric [98, 24].

Although many search space designs are possible and some have already been explored (Appendix B), there is an unexplored yet promising approach where we can define the entire agentic system in code and new agents can be automatically discovered by a meta agent programming even better ones in code. Searching within a code space theoretically enables the ADAS algorithm to discover *any* possible building blocks (e.g., prompts, tool use, workflow) and agentic systems that combine any of these building blocks in any way. This approach also offers better interpretability for agent design patterns since the program code is often readable, making debugging easier and enhancing AI safety. Additionally, compared to search spaces using networks [41] or graphs [98], searching in a code space allows us to more easily build on existing human efforts. For example, it is possible to search within open-source agent frameworks like LangChain [37] and build upon all existing building blocks (e.g., RAG, search engine tools). Finally, since FMs are proficient in coding, utilizing a code search space allows us to leverage existing expertise from FMs during the search process. In contrast, search algorithms in custom search spaces, such as graphs, may be much less efficient due to the absence of these priors. Therefore, we argue that the approach of using programming languages as the search space should be studied more in ADAS.

B Related Work

Agentic Systems. Researchers develop various building blocks and design patterns for different applications. Important building blocks for agentic systems includes: prompting techniques [8, 69], chain-of-thought-based planning and reasoning methods [83, 87, 30], reflection [47, 72], developing new skills for embodied agents in code [77, 76], external memory and RAG [94, 39], tool use [61, 68, 52], assigning FM modules in the agentic system with different roles and enabling them to collaborate [29, 84, 59, 85, 60], and enabling the agent to instruct itself for the next action [64], etc. While the community has invested substantial effort in developing all the above important

techniques, this is only a partial list of the discovered building blocks, and many more remain to be uncovered. Therefore, in this paper, we describe a newly forming research area, ADAS, which aims to invent novel building blocks and design powerful agentic systems in an automated manner.

AI-Generating Algorithms and AutoML. Research in AI-Generating Algorithms (AI-GAs) [14] and AutoML [33] aims to replace handcrafted components in AI systems by learning them. This field has three key pillars: (1) meta-learning architectures, (2) meta-learning learning algorithms, and (3) generating learning environments and training data [14]. Neural Architecture Search [22] exemplifies the first pillar by automating neural network design, while works like MAML [25] and Meta-RL [78] exemplify the second pillar, focusing on “learning to learn” for improved sample efficiency and generalizability. The third pillar includes works like POET [80] and OMNI-EPIC [23], which generate learning environments in an open-ended manner. We position Automated Design of Agentic Systems in both the first and second pillars: meta-learning agentic architectures and leveraging in-context learning to “learn to learn,” as shown in the ARC challenge (Section 3). Recent AI-GA and AutoML advances have also integrated Foundation Models (FMs) to write code, as seen in FunSearch [67] and EoH [40], where FMs discover optimization algorithms. In DiscoPOP [43], FMs program loss functions for preference learning, and Eureka [46] and language-to-reward [88] enable FMs to write reward functions for reinforcement learning. OMNI-EPIC [23] allows FMs to create robotics learning environments. Similarly, we enable FMs to program new agents in code.

Existing Attempts to ADAS. There are two categories of works that attempt ADAS: those focused on learning better prompts and those that learn more components beyond prompts. Most works fall into the first category, where FMs are used to automate prompt engineering, primarily enhancing the phrasing of instructions to improve reasoning [86, 24, 96]. However, these prompts are often domain-specific and difficult to generalize. Some works optimize role definitions within prompts [89, 11, 9, 84], but other components remain fixed, limiting the space of agents that can be discovered. The second category, which is less explored, involves learning additional components such as workflows, often representing agents as networks or graphs. For example, DyLAN [41] and DSPy [34] use FMs to optimize connections between nodes in a network, while GPT-Swarm [98] uses reinforcement learning to optimize node connections and prompts. Although these approaches optimize workflows, many components like tool usage or node structure remain fixed. AgentOptimizer [93] and Agent Symbolic Learning [97] attempt to learn prompts, tools, and workflows together, but often rely on predefined, complex agents. In contrast, our work allows all components to be represented in code, leveraging existing codebases and FMs’ expertise, which simplifies the search and enables the emergence of novel design patterns and building blocks from basic agent designs, showcasing ADAS’s creative potential.

C Future Work

Our work also paves the way for several future research directions:

- **Higher-order ADAS.** Since the meta agent used in ADAS to program new agents in code is also an agent, ADAS can become self-referential where the meta agent can be improved through ADAS as well. It would be an exciting direction to have a higher order of meta-learning to allow the learning of the meta agent and even the meta-meta agent, etc. [42]
- **Online Continual Learning.** As agents are deployed, they will receive vast amounts of feedback from both task environments and users. Continuously improving agents based on this extensive feedback is challenging for human developers. However, with ADAS automating the design and enhancement of agents, online continual learning becomes feasible post-deployment.
- **Multi-objective ADAS.** We only consider one objective (i.e., performance) to optimize in this paper, but in practice, multiple objectives are often considered, such as cost, latency, and robustness of agentic systems [31, 32]. Thus, integrating multi-objective search algorithms [18] in ADAS could be promising.
- **More complex domains.** Currently, we only evaluate Meta Agent Search on single-step QA tasks in this paper. It would be interesting to extend the method to more complex domains, such as real-world applications involving multi-step interaction with complex environments.
- **More intelligent evaluation functions.** In this work, we simply evaluate discovered agents on the evaluation set and use the numerical performance results. However, this approach is both

expensive and misses a lot of information. A promising future direction is to enable the meta agent to analyze detailed running logs during the evaluation, which contain rich information on the failure and success modes for better debugging and improving agentic systems [97]. Also, many tasks involve subjective answer evaluations [12, 44] that do not have ground-truth answers. It is also important to design novel evaluation functions in ADAS to address these tasks. Finally, in this work, we targeted only one domain during the search. It would be interesting to explore whether ADAS algorithms can design even better generalist agents when specifically searching for agents capable of performing well across multiple domains.

- **Towards a Better Understanding of FMs.** Works from Neural Architecture Search [32] show that by observing the emerged architecture, we could gain more insights into Neural Networks. In this paper, we also gained insights about FMs from the results. For example, the best agent with GPT-3.5 involves a complex feedback mechanism, but when we transfer to other advanced models, the agent with a simpler feedback mechanism but more refinement becomes a better agent (Section 3). This shows that GPT-3.5 may have a worse capability in evaluating and refining the answers, so it needs a complex feedback mechanism for better refinement, while other advanced models benefit more from a simpler feedback mechanism.
- **Seeding ADAS with more existing building blocks.** Although we can theoretically allow any components in agentic systems to be programmed from scratch in the code space, it is not efficient in practice. Therefore, it would be interesting to explore ADAS by standing on the shoulders of existing human efforts, such as search engine tools, RAG [39], or functions from existing agent frameworks like LangChain [37]. Additionally, it is interesting to support multi-modal capabilities (e.g. vision) in FMs or allow different FMs to be available in agentic systems. This will enable the meta agent to choose from different FMs flexibly according to the difficulty of the instruction and whether data privacy is a priority.
- **Novelty search algorithms.** In Meta Agent Search, the design of the search algorithm is relatively simple, focusing solely on exploring interesting new designs. A more careful design of the search algorithm can be a promising future direction. For example, one could incorporate more sophisticated ideas from Quality-Diversity [51, 16], AI-generating [14], and Open-ended Algorithms [23, 92, 73, 74]. One could also include more classic approaches to balance exploration and exploitation [75, 40].
- **Understanding the emergence of complexity from human organizations.** Beyond potentially saving researchers’ efforts and improving upon the manual design of agentic systems, the research in ADAS is also scientifically intriguing as it sheds light on the origins of complexity emerging from human organization and society. The agentic system is a machine learning system that operates primarily over natural language—a representation that is interpretable to humans and used by humans in constructing our organization and society. Thus, there is a close connection between agentic systems and human organizations, as shown in works incorporating the organizational structure for human companies in agents [29] or simulating a human town with agents [57]. Therefore, the study in ADAS may enable us to observe how to create a simple set of conditions and have an algorithm to bootstrap itself from simplicity to produce complexity in a system akin to human society.

D Generalizability and Transferability

In this section, we present the complete results of transferring agents across different models and domains.

Transferability Across Foundation Models. We first transfer discovered agents from GPT-3.5 [54] to other FMs on ARC to test whether agents found when performing Meta Agent Search with one FM generalize to others. We test the top 3 agents with the best test accuracy evaluated with GPT-3.5 on ARC and then transfer them to three popular models: Claude-Haiku [1], GPT-4 [56], and Claude-Sonnet [2]. We adopt the same baselines as those used in ARC (Section 3) and MGSM (Section 3). As shown in Table 3, we observe that the searched agents consistently outperform the hand-designed agents with a substantial gap. Notably, we found that Claude-Sonnet, the most powerful model from Anthropic, performs the best among all tested models, enabling our best agent to achieve nearly 50% accuracy on ARC.

Table 3: **Performance on ARC when transferring top agents from GPT-3.5 to other FMs.** Agents discovered by Meta Agent Search consistently outperform the baselines across different models. We report the test accuracy and the 95% bootstrap confidence interval. The names of top agents are generated by Meta Agent Search. †We manually changed this name because the original generated name was confusing.

Agent Name	Accuracy on ARC (%)			
	GPT-3.5	Claude-Haiku	GPT-4	Claude-Sonnet
Manually Designed Agents				
Chain-of-Thought [83]	6.0 ± 2.7	4.3 ± 2.2	17.7 ± 4.4	25.3 ± 5.0
COT-SC [82]	8.0 ± 3.2	5.3 ± 2.5	19.7 ± 4.5	26.3 ± 4.9
LLM Debate [19]	4.0 ± 2.2	1.7 ± 1.5	19.0 ± 4.5	24.7 ± 4.8
Self-Refine [47]	6.7 ± 2.7	6.3 ± 2.8	23.0 ± 5.2	39.3 ± 5.5
Quality-Diversity [45]	7.0 ± 2.9	3.3 ± 2.2	23.0 ± 4.7	31.7 ± 5.3
Top Agents Searched with GPT-3.5		Transferred to Other FMs		
Structured Feedback and Ensemble Agent	13.7 ± 3.9	5.0 ± 2.5	30.0 ± 5.2	38.7 ± 5.5
Hierarchical Committee Reinforcement Agent	13.3 ± 3.8	8.3 ± 3.2	32.3 ± 8.9	39.7 ± 5.5
Dynamic Memory and Refinement Agent†	12.7 ± 3.9	9.7 ± 3.3	37.0 ± 5.3	48.3 ± 5.7

Table 4: **Performance on different math domains when transferring top agents from MGSM to other math domains.** Agents discovered by Meta Agent Search consistently outperform the baselines across different math domains. We report the test accuracy and the 95% bootstrap confidence interval. The names of top agents are generated by Meta Agent Search.

Agent Name	Accuracy (%)				
	MGSM	GSM8K	GSM-Hard	SVAMP	ASDiv
Manually Designed Agents					
Chain-of-Thought [83]	28.0 ± 3.1	34.9 ± 3.2	15.0 ± 2.5	77.8 ± 2.8	88.9 ± 2.2
COT-SC [82]	28.2 ± 3.1	37.8 ± 3.4	15.5 ± 2.5	78.2 ± 2.8	89.0 ± 2.1
Self-Refine [47]	27.5 ± 3.1	38.9 ± 3.4	15.1 ± 2.4	78.5 ± 2.8	89.2 ± 2.2
LLM Debate [19]	39.0 ± 3.4	43.6 ± 3.4	17.4 ± 2.6	76.0 ± 3.0	88.9 ± 2.2
Step-back Abstraction [95]	31.1 ± 3.2	31.5 ± 3.3	12.2 ± 2.3	76.1 ± 3.0	87.8 ± 2.3
Quality-Diversity [45]	23.8 ± 3.0	28.0 ± 3.1	14.1 ± 2.4	69.8 ± 3.2	80.1 ± 2.8
Role Assignment [85]	30.1 ± 3.2	37.0 ± 3.4	18.0 ± 2.7	73.0 ± 3.0	83.1 ± 2.6
Top Agents Searched on MGSM (Math)		Transferred within Math Domains			
Dynamic Role-Playing Architecture	53.4 ± 3.5	69.5 ± 3.2	31.2 ± 3.2	81.5 ± 2.6	91.8 ± 1.8
Structured Multimodal Feedback Loop	50.2 ± 3.5	64.5 ± 3.4	30.1 ± 3.2	82.6 ± 2.6	89.9 ± 2.1
Interactive Multimodal Feedback Loop	47.4 ± 3.5	64.9 ± 3.3	27.6 ± 3.2	80.6 ± 2.8	89.8 ± 2.1

Transferability Across Domains. Next, we transfer the discovered agent from the MGSM (Math) domain to other math domains to test whether the invented agents can generalize across different domains. Similarly, we test the top 3 agents from MGSM and transfer them to (1) four popular math domains: GSM8K [15], GSM-Hard [26], SVAMP [58], and ASDiv [50] and (2) three domains beyond math adopted in Section 3. As shown in Table 4, we observe a similar superiority in the performance of Meta Agent Search compared to baselines. More surprisingly, we observe that agents discovered in the math domain can be transferred to non-math domains (Table 5). While the performance of agents originally searched in the math domain does not fully match that of agents specifically designed for the target domains, they still outperform (in Reading Comprehension and Multi-task) or match (in Science) the state-of-the-art hand-designed agent baselines. These results illustrate that Meta Agent Search can discover generalizable design patterns and agentic systems.

E Prompts

We use the following prompts for the meta agent in Meta Agent Search. Variables in the prompts that vary depending on domains and iterations are **highlighted**.

We use the following system prompt for every query in the meta agent.

Table 5: **Performance across multiple domains when transferring top agents from the Math (MGSM) domain to non-math domains.** Agents discovered by Meta Agent Search in the math domain can outperform or match the performance of baselines after being transferred to domains beyond math. We report the test accuracy and the 95% bootstrap confidence interval.

Agent Name	Accuracy (%)		F1 Score	
	Math	Reading Comprehension	Multi-task	Science
Manually Designed Agents				
Chain-of-Thought [83]	28.0 ± 3.1	64.2 ± 0.9	65.4 ± 3.3	29.2 ± 3.1
COT-SC [82]	28.2 ± 3.1	64.4 ± 0.8	65.9 ± 3.2	30.5 ± 3.2
Self-Refine [47]	27.5 ± 3.1	59.2 ± 0.9	63.5 ± 3.4	31.6 ± 3.2
LLM Debate [19]	39.0 ± 3.4	60.6 ± 0.9	65.6 ± 3.3	31.4 ± 3.2
Step-back Abstraction [95]	31.1 ± 3.2	60.4 ± 1.0	65.1 ± 3.3	26.9 ± 3.0
Quality-Diversity [45]	23.8 ± 3.0	61.8 ± 0.9	65.1 ± 3.1	30.2 ± 3.1
Role Assignment [85]	30.1 ± 3.2	65.8 ± 0.9	64.5 ± 3.3	31.1 ± 3.1
Top Agents Searched on Math (MGSM)		Transferred beyond Math Domains		
Dynamic Role-Playing Architecture	53.4 ± 3.5	70.4 ± 0.9	62.4 ± 3.4	28.6 ± 3.1
Structured Multimodal Feedback Loop	50.2 ± 3.5	70.4 ± 0.9	67.0 ± 3.2	28.7 ± 3.1
Interactive Multimodal Feedback Loop	47.4 ± 3.5	71.9 ± 0.8	64.8 ± 3.3	29.9 ± 3.2

System prompt for the meta agent.

You are a helpful assistant. Make sure to return in a WELL-FORMED JSON object.

We use the following prompt for the meta agent to design the new agent based on the archive of previously discovered agents.

Main prompt for the meta agent.

You are an expert machine learning researcher testing various agentic systems. Your objective is to design building blocks such as prompts and workflows within these systems to solve complex tasks. Your aim is to design an optimal agent performing well on **[Brief Description of the Domain]**.

[Framework Code]

[Output Instructions and Examples]

[Discovered Agent Archive] (initialized with baselines, updated at every iteration)

Your task

You are deeply familiar with prompting techniques and the agent works from the literature. Your goal is to maximize the specified performance metrics by proposing interestingly new agents.

Observe the discovered agents carefully and think about what insights, lessons, or stepping stones can be learned from them.

Be creative when thinking about the next interesting agent to try. You are encouraged to draw inspiration from related agent papers or academic papers from other research areas.

Use the knowledge from the archive and inspiration from academic literature to propose the next interesting agentic system design.

THINK OUTSIDE THE BOX.

The domain descriptions are available in Appendices G and H and the framework code is available in Appendix F. We use the following prompt to instruct and format the output of the meta agent. Here, we collect and present some common mistakes that the meta agent may make in the prompt. We found it effective in improving the quality of the generated code.

Output Instruction and Example.

Output Instruction and Example:

The first key should be (“thought”), and it should capture your thought process for designing the next function. In the “thought” section, first reason about what the next interesting agent to try should be, then describe your reasoning and the overall concept behind the agent design, and finally detail

the implementation steps. The second key (“name”) corresponds to the name of your next agent architecture. Finally, the last key (“code”) corresponds to the exact “forward()” function in Python code that you would like to try. You must write COMPLETE CODE in “code”: Your code will be part of the entire project, so please implement complete, reliable, reusable code snippets.

Here is an example of the output format for the next agent:

```
{“thought”: “**Insights:** Your insights on what should be the next interesting agent. **Overall Idea:** your reasoning and the overall concept behind the agent design. **Implementation:** describe the implementation step by step.”,  
“name”: “Name of your proposed agent”,  
“code”: “def forward(self, taskInfo): # Your code here”}
```

WRONG Implementation examples:

[Examples of potential mistakes the meta agent may make in implementation]

After the first response from the meta agent, we perform two rounds of self-reflection to make the generated agent novel and error-free [72, 47].

Prompt for self-reflection round 1.

[Generated Agent from Previous Iteration]

Carefully review the proposed new architecture and reflect on the following points:

- Interestingness**: Assess whether your proposed architecture is interesting or innovative compared to existing methods in the archive. If you determine that the proposed architecture is not interesting, suggest a new architecture that addresses these shortcomings.
 - Make sure to check the difference between the proposed architecture and previous attempts.
 - Compare the proposal and the architectures in the archive CAREFULLY, including their actual differences in the implementation.
 - Decide whether the current architecture is innovative.
 - USE CRITICAL THINKING!
- Implementation Mistakes**: Identify any mistakes you may have made in the implementation. Review the code carefully, debug any issues you find, and provide a corrected version. REMEMBER checking “## WRONG Implementation examples” in the prompt.
- Improvement**: Based on the proposed architecture, suggest improvements in the detailed implementation that could increase its performance or effectiveness. In this step, focus on refining and optimizing the existing implementation without altering the overall design framework, except if you want to propose a different architecture if the current is not interesting.
 - Observe carefully about whether the implementation is actually doing what it is supposed to do.
 - Check if there is redundant code or unnecessary steps in the implementation. Replace them with effective implementation.
 - Try to avoid the implementation being too similar to the previous agent.

And then, you need to improve or revise the implementation, or implement the new proposed architecture based on the reflection.

Your response should be organized as follows:

”reflection”: Provide your thoughts on the interestingness of the architecture, identify any mistakes in the implementation, and suggest improvements.

”thought”: Revise your previous proposal or propose a new architecture if necessary, using the same format as the example response.

”name”: Provide a name for the revised or new architecture. (Don’t put words like ”new” or ”improved” in the name.)

”code”: Provide the corrected code or an improved implementation. Make sure you actually implement your fix and improvement in this code.

Prompt for self-reflection round 2.

Using the tips in “## WRONG Implementation examples” section, further revise the code. Your response should be organized as follows:
Include your updated reflections in the “reflection”. Repeat the previous “thought” and “name”. Update the corrected version of the code in the “code” section.

When an error is encountered during the execution of the generated code, we conduct a reflection and re-run the code. This process is repeated up to five times if errors persist. Here is the prompt we use to self-reflect any runtime error:

Prompt for self-reflection when a runtime error occurs.

Error during evaluation:

[Runtime errors]

Carefully consider where you went wrong in your latest implementation. Using insights from previous attempts, try to debug the current code to implement the same thought. Repeat your previous thought in “thought”, and put your thinking for debugging in “debug.thought”.

F Framework Code

In this paper, we provide the meta agent with a simple framework to implement basic functions, such as querying Foundation Models (FMs) and formatting prompts. The framework consists of fewer than 100 lines of code (excluding comments). In this framework, we encapsulate every piece of information into a namedtuple Info object, making it easy to combine different types of information (e.g., FM responses, results from tool function calls, task descriptions) and facilitate communication between different modules. Additionally, in the FM module, we automatically construct the prompt by concatenating all input Info objects into a structured format, with each Info titled by its metadata (e.g., name, author). **Throughout the appendix, we renamed some variables in the code to match the terminologies used in the main text.**

Code 1: The simple framework used in Meta-Agent Search.

```
1 # Named tuple for holding task information
2 Info = namedtuple('Info', ['name', 'author', 'content', 'iteration_idx',
3                             ''])
4
5 # Format instructions for FM response
6 FORMAT_INST = lambda request_keys: f"Reply EXACTLY with the following
7     JSON format.\n{str(request_keys)}\nDO NOT MISS ANY FIELDS AND MAKE
8     SURE THE JSON FORMAT IS CORRECT!\n"
9
10 # Description of the role of the FM Module
11 ROLE_DESC = lambda role: f"You are a {role}."
12
13 @backoff.on_exception(backoff.expo, openai.RateLimitError)
14 def get_json_response_from_gpt(msg, model, system_message, temperature):
15     \"""
16     Function to get JSON response from GPT model.
17
18     Args:
19     - msg (str): The user message.
20     - model (str): The model to use.
21     - system_message (str): The system message.
22     - temperature (float): Sampling temperature.
23
24     Returns:
25     - dict: The JSON response.
26     \"""
27     ...
28     return json_dict
```

```

27 class FM_Module:
28     \"""
29     Base class for an FM module.
30
31     Attributes:
32     - output_fields (list): Fields expected in the output.
33     - name (str): Name of the FM module.
34     - role (str): Role description for the FM module.
35     - model (str): Model to be used.
36     - temperature (float): Sampling temperature.
37     - id (str): Unique identifier for the FM module instance.
38     \"""
39
40     def __init__(self, output_fields: list, name: str, role='helpful
41         assistant', model='gpt-3.5-turbo-0125', temperature=0.5) ->
42         None:
43         ...
44
45     def generate_prompt(self, input_infos, instruction) -> str:
46         \"""
47         Generates a prompt for the FM.
48
49         Args:
50         - input_infos (list): List of input information.
51         - instruction (str): Instruction for the task.
52
53         Returns:
54         - tuple: System prompt and user prompt.
55
56         An example of generated prompt:
57         ""
58         You are a helpful assistant.
59
60         # Output Format:
61         Reply EXACTLY with the following JSON format.
62         ...
63
64         # Your Task:
65         You will given some number of paired example inputs and
66         outputs. The outputs ...
67
68         ### thinking #1 by Chain-of-Thought hkFo (yourself):
69         ...
70
71         # Instruction:
72         Please think step by step and then solve the task by writing
73         the code.
74         ""
75         \"""
76         ...
77         return system_prompt, prompt
78
79     def query(self, input_infos: list, instruction, iteration_idx=-1)
80         -> list[Info]:
81         \"""
82         Queries the FM with provided input information and instruction
83         .
84
85         Args:
86         - input_infos (list): List of input information.
87         - instruction (str): Instruction for the task.
88         - iteration_idx (int): Iteration index for the task.
89
90         Returns:
91         - output_infos (list[Info]): Output information.

```

```

86     \"""
87     ...
88     return output_infos
89
90     def __repr__(self):
91         return f"{self.agent_name} {self.id}"
92
93     def __call__(self, input_infos: list, instruction, iteration_idx
94         =-1):
95         return self.query(input_infos, instruction, iteration_idx=
96             iteration_idx)
97
98 class AgentSystem:
99     def forward(self, taskInfo) -> Union[Info, str]:
100         \"""
101         Placeholder method for processing task information.
102
103         Args:
104         - taskInfo (Info): Task information.
105
106         Returns:
107         - Answer (Union[Info, str]): Your FINAL Answer. Return either
108             a namedtuple Info or a string for the answer.
109         \"""
110         pass

```

With the provided framework, an agent can be easily defined with a “forward” function. Here we show an example of implementing self-reflection using the framework.

Code 2: Self-Reflection implementation example

```

1  def forward(self, taskInfo):
2      # Instruction for initial reasoning
3      cot_initial_instruction = "Please think step by step and then
4          solve the task."
5
6      # Instruction for reflecting on previous attempts and feedback to
7          improve
8      cot_reflect_instruction = "Given previous attempts and feedback,
9          carefully consider where you could go wrong in your latest
10         attempt. Using insights from previous attempts, try to solve
11         the task better."
12      cot_module = FM_Module(['thinking', 'answer'], 'Chain-of-Thought')
13
14      # Instruction for providing feedback and correcting the answer
15      critic_instruction = "Please review the answer above and criticize
16         on where might be wrong. If you are absolutely sure it is
17         correct, output 'True' in 'correct'."
18      critic_module = FM_Module(['feedback', 'correct'], 'Critic')
19
20      N_max = 5 # Maximum number of attempts
21
22      # Initial attempt
23      cot_inputs = [taskInfo]
24      thinking, answer = cot_module(cot_inputs, cot_initial_instruction,
25          0)
26
27      for i in range(N_max):
28          # Get feedback and correct status from the critic
29          feedback, correct = critic_module([taskInfo, thinking, answer
30              ], critic_instruction, i)
31          if correct.content == 'True':
32              break
33
34          # Add feedback to the inputs for the next iteration

```

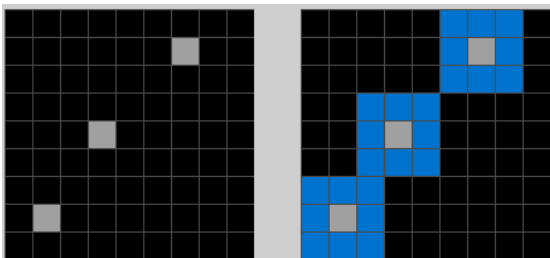
```

26     cot_inputs.extend([thinking, answer, feedback])
27
28     # Reflect on previous attempts and refine the answer
29     thinking, answer = cot_module(cot_inputs,
30     cot_reflect_instruction, i + 1)
    return answer

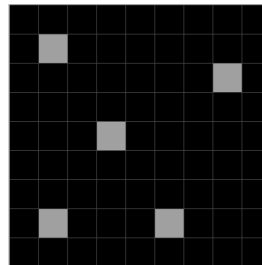
```

G Experiment Details for ARC Challenge

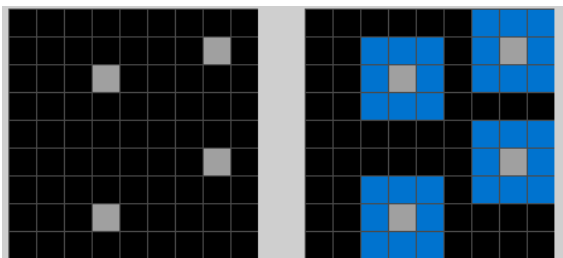
Example Input-output grid #1



Test grid



Example Input-output grid #2



Answer

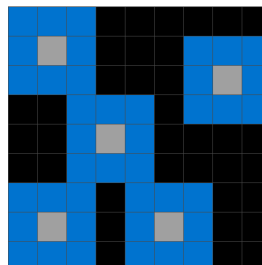


Figure 4: **An example task from the ARC challenge [13].** Given the input-output grid examples, the AI system is asked to learn the transformation rules and then apply these learned rules to the test grid to predict the final answer.

The ARC challenge aims to evaluate the general intelligence of AI systems through their ability to acquire new skills. Questions in ARC include (1) showing multiple examples of visual input-output grid patterns, (2) the AI system learning the transformation rule of grid patterns from examples, and (3) predicting the output grid pattern given a test input grid pattern. Since each question in ARC has a unique transformation rule, it requires the AI system to learn efficiently with few-shot examples, leveraging capabilities in number counting, geometry, and topology. Following common practice [27], we require the agent to write code for the transformation rule instead of answering directly. We provide tool functions in the framework that evaluate the generated transformation code. Given the significant challenge that ARC poses to current AI systems, we sample our data from questions with grid dimensions $\leq 5 \times 5$ in the “Public Training Set (Easy)”. We sample a validation set and a test set with 20 and 60 questions, respectively, for searching and testing. We calculate the validation and test accuracy of an agent by assessing it over the validation and test sets five times to reduce the variance from the stochastic sampling of FMs. We evaluate all discovered agents on the held-out test set and report the test accuracy in Figure 2. Meta Agent Search runs for 25 iterations and the meta agent uses GPT-4 [56], while discovered agents and baselines are evaluated using GPT-3.5 [54] to reduce compute cost. We compare against five state-of-the-art hand-designed agents, and we also use all baselines as initial seeds in the archive for Meta Agent Search.

An example task from the ARC challenge is shown in Figure 4. In the ARC challenge experiments (Section 3), we represent the grids as strings of 2-D arrays, where each color is represented by an

integer. We instruct the meta agent to design agents that generate code as solutions rather than directly outputting answers. Additionally, we provide two tool functions within the framework: (1) to test whether the generated code can solve the example grids and (2) to obtain the task’s answer by applying the generated code to the test grid. The accuracy rate is calculated by the Exact Match between the reference solution and the predicted answer. The meta agent uses “gpt-4o-2024-05-13” [56], while discovered agents and baselines are evaluated using “gpt-3.5-turbo-0125” [54] to reduce compute cost.

The domain description of ARC for the meta agent is shown below:

```

Description of ARC for the meta agent.

Your aim is to find an optimal agent performing well on the ARC (Abstraction and Reasoning Corpus) challenge.
In this challenge, each task consists of three demonstration examples, and one test example. Each Example consists of an “input grid” and an “output grid”. Test-takers need to use the transformation rule learned from the examples to predict the output grid for the test example.

# An example task from ARC challenge:

## Task Overview:
You will be given some number of paired example inputs and outputs grids. The outputs were produced by applying a transformation rule to the input grids. In addition to the paired example inputs and outputs, there is also one test input without a known output.
The inputs and outputs are each “grids”. A grid is a rectangular matrix of integers between 0 and 9 (inclusive). Each number corresponds to a color. 0 is black.
Your task is to determine the transformation rule from examples and find out the answer, involving determining the size of the output grid for the test and correctly filling each cell of the grid with the appropriate color or number.

The transformation only needs to be unambiguous and applicable to the example inputs and the test input. It doesn’t need to work for all possible inputs. Observe the examples carefully, imagine the grid visually, and try to find the pattern.

## Examples:
### Example 0:
input = [[0,0,0,0,5,0,0,0,0], [0,0,0,0,5,0,0,0,0], [0,0,0,4,5,0,0,0,0], [0,0,0,4,5,4,4,0,0], [0,0,3,3,5,0,0,0,0],
[0,0,0,3,5,0,0,0,0], [0,0,0,3,5,3,3,3,0], [0,0,0,3,5,0,0,0,0], [0,0,0,0,5,0,0,0,0], [0,0,0,0,5,0,0,0,0]]
output = [[0,0,0,0], [0,0,0,0], [0,0,0,4], [0,0,4,4], [0,0,3,3], [0,0,0,3], [0,3,3,3], [0,0,0,3], [0,0,0,0],
[0,0,0,0]]

### Example 1:
input = [[0,0,0,0,5,0,0,0,0], [0,0,0,2,5,0,0,0,0], [0,0,0,2,5,2,6,0,0], [0,0,0,2,5,0,0,0,0], [0,0,0,2,5,2,2,2,0],
[0,0,6,6,5,6,0,0,0], [0,0,0,2,5,0,0,0,0], [0,2,2,0,5,2,0,0,0], [0,0,0,2,5,0,0,0,0], [0,0,0,0,5,0,0,0,0]]
output = [[0,0,0,0], [0,0,0,2], [0,0,6,2], [0,0,0,2], [0,2,2,2], [0,0,6,6], [0,0,0,2], [0,2,2,2], [0,0,0,2],
[0,0,0,0]]

### Example 2:
input = [[0,0,0,0,5,0,0,0,0], [0,0,0,0,5,7,0,0,0], [0,0,0,8,5,0,0,0,0], [0,0,0,8,5,0,0,0,0], [0,7,8,8,5,0,0,0,0],
[0,0,0,0,5,8,8,0,0], [0,0,0,8,5,0,0,0,0], [0,0,0,8,5,0,0,0,0], [0,0,0,0,5,8,7,0,0], [0,0,0,0,5,0,0,0,0]]
output= [[0,0,0,0], [0,0,0,7], [0,0,0,8], [0,0,0,8], [0,7,8,8], [0,0,8,8], [0,0,0,8], [0,0,0,8], [0,0,7,8],
[0,0,0,0]]

### Test Problem:
input = [[0,0,0,0,5,0,0,0,0], [0,0,0,1,5,0,0,0,0], [0,0,0,1,5,1,0,0,0], [0,1,1,1,5,1,1,1,6], [0,0,0,6,5,6,6,0,0],
[0,0,0,0,5,1,1,1,0], [0,0,0,1,5,0,0,0,0], [0,0,0,1,5,1,6,0,0], [0,0,0,0,5,6,0,0,0], [0,0,0,0,5,0,0,0,0]]

Analyze the transformation rules based on the provided Examples and determine what the output should be for the Test Problem.

```

Here we present the best agent on ARC discovered by Meta Agent Search.

Code 3: The best agent on ARC discovered by Meta Agent Search

```

1 # Structured Feedback and Ensemble Agent
2 def forward(self, taskInfo):
3     # Step 1: Generate initial candidate solutions using multiple FM
4     # Modules
5     initial_instruction = 'Please think step by step and then solve
6     the task by writing the code.'
7     num_candidates = 5 # Number of initial candidates
8     initial_module = [FM_Module(['thinking', 'code'], 'Initial
9     Solution', temperature=0.8) for _ in range(num_candidates)]
10
11     initial_solutions = []
12     for i in range(num_candidates):
13         thoughts = initial_module[i]([taskInfo], initial_instruction)
14         thinking, code = thoughts[0], thoughts[1]
15         feedback, correct_examples, wrong_examples = self.
16         run_examples_and_get_feedback(code)
17         if len(correct_examples) > 0: # Only consider solutions that
18         passed at least one example
19             initial_solutions.append({'thinking': thinking, 'code':
20             code, 'feedback': feedback, 'correct_count': len(
21             correct_examples)})
22
23     # Step 2: Simulate human-like feedback for each candidate solution
24     human_like_feedback_module = FM_Module(['thinking', 'feedback'], '
25     Human-like Feedback', temperature=0.5)
26     human_feedback_instruction = 'Please provide human-like feedback
27     for the code, focusing on common mistakes, heuristic
28     corrections, and best practices.'
29
30     for sol in initial_solutions:
31         thoughts = human_like_feedback_module([taskInfo, sol['thinking
32         '], sol['code']], human_feedback_instruction)
33         human_thinking, human_feedback = thoughts[0], thoughts[1]
34         sol['human_feedback'] = human_feedback
35
36     # Step 3: Assign expert advisors to evaluate and provide targeted
37     # feedback
38     expert_roles = ['Efficiency Expert', 'Readability Expert', '
39     Simplicity Expert']
40     expert_advisors = [FM_Module(['thinking', 'feedback'], role,
41     temperature=0.6) for role in expert_roles]
42     expert_instruction = 'Please evaluate the given code and provide
43     targeted feedback for improvement.'
44
45     for sol in initial_solutions:
46         sol_feedback = {}
47         for advisor in expert_advisors:
48             thoughts = advisor([taskInfo, sol['thinking'], sol['code'
49             ']], expert_instruction)
50             thinking, feedback = thoughts[0], thoughts[1]
51             sol_feedback[advisor.role] = feedback
52         sol['expert_feedback'] = sol_feedback
53
54     # Step 4: Parse and structure the feedback to avoid redundancy and
55     # refine the solutions iteratively
56     max_refinement_iterations = 3
57     refinement_module = FM_Module(['thinking', 'code'], 'Refinement
58     Module', temperature=0.5)
59     refined_solutions = []
60
61     for sol in initial_solutions:
62         for i in range(max_refinement_iterations):
63             combined_feedback = sol['feedback'].content + sol['
64             human_feedback'].content + ''.join([fb.content for fb
65             in sol['expert_feedback'].values()])

```



```

46     structured_feedback = ' '.join(set(combined_feedback.split
47         ())) # Avoid redundancy
48     refinement_instruction = 'Using the structured feedback,
49         refine the solution to improve its performance.'
50     thoughts = refinement_module([taskInfo, sol['thinking'],
51         sol['code'], Info('feedback', 'Structured Feedback',
52         structured_feedback, i)], refinement_instruction, i)
53     refinement_thinking, refined_code = thoughts[0], thoughts
54         [1]
55     feedback, correct_examples, wrong_examples = self.
56         run_examples_and_get_feedback(refined_code)
57     if len(correct_examples) > 0:
58         sol.update({'thinking': refinement_thinking, 'code':
59             refined_code, 'feedback': feedback, 'correct_count
60             ': len(correct_examples)})
61         refined_solutions.append(sol)
62
63     # Step 5: Select the best-performing solutions and make a final
64     # decision using an ensemble approach
65     sorted_solutions = sorted(refined_solutions, key=lambda x: x['
66         correct_count'], reverse=True)
67     top_solutions = sorted_solutions[:3] # Select the top 3 solutions
68
69     final_decision_instruction = 'Given all the above solutions,
70         reason over them carefully and provide a final answer by
71         writing the code.'
72     final_decision_module = refinement_module(['thinking', 'code'], '
73         Final Decision Module', temperature=0.1)
74     final_inputs = [taskInfo] + [item for solution in top_solutions
75         for item in [solution['thinking'], solution['code'], solution[
76             'feedback']]]
77     final_thoughts = final_decision_module(final_inputs,
78         final_decision_instruction)
79     final_thinking, final_code = final_thoughts[0], final_thoughts[1]
80     answer = self.get_test_output_from_code(final_code)
81     return answer

```

H Experiment Details for Reasoning and Problem-Solving Domains

We test Meta Agent Search on four popular benchmarks: (1) DROP [20] for evaluating **Reading Comprehension**; (2) MGSM [71] for evaluating **Math** capability under a multi-lingual setting; (3) MMLU [28] for evaluating **Multi-task** Problem Solving; and (4) GPQA [63] for evaluating the capability of solving hard (graduate-level) questions in **Science**. The search is conducted independently within each domain. Meta Agent Search runs for 30 iterations. The meta agent uses GPT-4 [56], while the discovered agents and baselines are evaluated using GPT-3.5 [54]. We adopt all baselines introduced in Section 3. Additionally, since the above domains require strong reasoning skills, we include two additional baselines Step-back Abstraction [95] and Role Assignment [85] that specifically focus on enhancing the reasoning capabilities of agents for a more thorough comparison.

To reduce costs during search and evaluation, we sample subsets of data from each domain. For GPQA (Science), we use GPQA_diamond and the validation set consists of 32 questions, while the remaining 166 questions form the test set. For the other domains, the validation and test sets are sampled with 128 and 800 questions, respectively. We evaluate agents five times for GPQA and once for the other domains to maintain a consistent total number of evaluations. Each domain uses zero-shot style questions, except DROP (Reading Comprehension), which uses one-shot style questions following the practice in [55]. The meta agent uses “gpt-4o-2024-05-13” [56], while discovered agents and baselines are evaluated using “gpt-3.5-turbo-0125” [54] to reduce compute cost.

We present the description of each domain we provide to the meta agent.

Description of DROP (Reading Comprehension).

Your aim is to find an optimal agent performing well on the Reading Comprehension Benchmark Requiring Discrete Reasoning Over Paragraphs (DROP), which assesses the ability to perform discrete reasoning and comprehend detailed information across multiple paragraphs.

An example question from DROP:

You will be asked to read a passage and answer a question.

Passage:

Non-nationals make up more than half of the population of Bahrain, with immigrants making up about 55% of the overall population. Of those, the vast majority come from South and Southeast Asia: according to various media reports and government statistics dated between 2005-2009 roughly 290,000 Indians, 125,000 Bangladeshis, 45,000 Pakistanis, 45,000 Filipinos, and 8,000 Indonesians.

Question: What two nationalities had the same number of people living in Bahrain between 2005-2009?

Answer [Not Given]: Pakistanis and Filipinos

Description of GPQA (Science) for the meta agent.

Your aim is to find an optimal agent performing well on the GPQA (Graduate-Level Google-Proof Q&A Benchmark). This benchmark consists of challenging multiple-choice questions across the domains of biology, physics, and chemistry, designed by domain experts to ensure high quality and difficulty.

An example question from GPQA:

Two quantum states with energies E_1 and E_2 have a lifetime of 10^{-9} sec and 10^{-8} sec, respectively. We want to clearly distinguish these two energy levels. Which one of the following options could be their energy difference so that they be clearly resolved?

Answer choices:

10^{-9} eV

10^{-8} eV

10^{-7} eV

10^{-6} eV

Correct answer [Not provided]:

10^{-7} eV

Explanation [Not provided]:

According to the uncertainty principle, $\Delta E \cdot \Delta t = \hbar/2$. Δt is the lifetime and ΔE is the width of the energy level. With $\Delta t = 10^{-9}$ s $\implies \Delta E_1 = 3.3 \cdot 10^{-7}$ eV. And $\Delta t = 10^{-8}$ s gives $\Delta E_2 = 3.3 \cdot 10^{-8}$ eV. Therefore, the energy difference between the two states must be significantly greater than 10^{-7} eV. So the answer is 10^{-4} eV.

Description of MGSM (Math) for the meta agent.

Your aim is to find an optimal agent performing well on the Multilingual Grade School Math Benchmark (MGSM) which evaluates mathematical problem-solving abilities across various languages to ensure broad and effective multilingual performance.

An example question from MGSM:

Question: この数学の問題を解いてください。

近所では、ペットのウサギの数がペットの犬と猫を合わせた数よりも12匹少ない。犬1匹あたり2匹の猫がおり、犬の数は60匹だとすると、全部で近所には何匹のペットがいますか？

Answer (Not Given): 348

Description of MMLU (Multi-task) for the meta agent.

Your aim is to find an optimal agent performing well on the MMLU (Massive Multitask Language Understanding) benchmark, a challenging evaluation that assesses a model’s ability to answer questions across a wide range of subjects and difficulty levels. It includes subjects from STEM, social sciences, humanities, and more.

An example question from MMLU:

Answer the following multiple-choice question.

The constellation ... is a bright W-shaped constellation in the northern sky.

- (A) Centaurus
- (B) Cygnus
- (C) Cassiopeia
- (D) Cepheus

I Baselines

In this paper, we implement five state-of-the-art hand-designed agent baselines for experiments on ARC (Section 3): (1) Chain-of-Thought (COT) [83], which instructs the agent to output the reasoning before answering to improve complex problem-solving through intermediate steps; (2) Self-Consistency with Chain-of-Thought (COT-SC) [82], which ensembles multiple parallel answers from COT to produce a more accurate answer; (3) Self-Refine [47, 72], which allows iterative self-reflection to correct mistakes made in previous attempts; (4) LLM-Debate [19], which enables different LLMs to debate with each other, leveraging diverse perspectives to find better answers; (5) Quality-Diversity, a simplified version of Intelligent Go-Explore [45], which produces and ensembles diverse answers to better explore potential solutions.

In addition to these baselines, we implement two more for experiments on Reasoning and Problem-Solving domains (Section 3): (6) Step-back Abstraction [95], which instructs agents to first consider the principles involved in solving the task for better reasoning and (7) Role Assignment [85], which assigns different roles to FMs to obtain better answers.. An example implementation of Self-Refine with our simple framework is shown in Appendix F.

In COT, we prompt the FM to think step by step before answering the question. In COT-SC, we sample $N = 5$ answers and then perform an ensemble using either majority voting or an FM query. In Self-Refine, we allow up to five refinement iterations, with an early stop if the critic deems the answer correct. In LLM-Debate, each debate module is assigned a unique role, such as Physics Expert or Chemistry Expert, and the debate lasts for two rounds. In Quality-Diversity, we conduct three iterations to collect diverse answers based on previously proposed ones. In Role Assignment, we use an FM query to first choose a role from a predefined set, and then use another FM query to answer the question by acting within the chosen role.

J Example Agents

In this section, we present the detailed implementation of three example discovered agents by Meta Agent Search shown in Figure 1. The “Multi-Step Peer Review Agent” and “Divide and Conquer Agent” were discovered during the search in the Reading Comprehension domain (GPQA) [63], while the “Verified Multimodal Agent” was discovered during the search in the Math domain (MGSM) [71].

Code 4: Example discovered agent: Multi-Step Peer Review Agent

```
1 def forward(self, taskInfo):  
2     initial_instruction = "Please think step by step and then solve  
    the task."
```

```

3 critique_instruction = "Please review the answer above and provide
  feedback on where it might be wrong. If you are absolutely
  sure it is correct, output 'True' in 'correct'."
4 refine_instruction = "Given previous attempts and feedback,
  carefully consider where you could go wrong in your latest
  attempt. Using insights from previous attempts, try to solve
  the task better."
5 final_decision_instruction = "Given all the above thinking and
  answers, reason over them carefully and provide a final answer
  ."
6
7 FM_modules = [FM_module(['thinking', 'answer'], 'FM Module', role=
  role) for role in ['Physics Expert', 'Chemistry Expert', '
  Biology Expert', 'Science Generalist']]
8 critic_modules = [FM_module(['feedback', 'correct'], 'Critic',
  role=role) for role in ['Physics Critic', 'Chemistry Critic',
  'Biology Critic', 'General Critic']]
9 final_decision_module = FM_module(['thinking', 'answer'], 'Final
  Decision', temperature=0.1)
10
11 all_thinking = [[] for _ in range(len(FM_modules))]
12 all_answer = [[] for _ in range(len(FM_modules))]
13 all_feedback = [[] for _ in range(len(FM_modules))]
14
15 for i in range(len(FM_modules)):
16     thinking, answer = FM_modules[i]([taskInfo],
17     initial_instruction)
18     all_thinking[i].append(thinking)
19     all_answer[i].append(answer)
20
21 for i in range(len(FM_modules)):
22     for j in range(len(FM_modules)):
23         if i != j:
24             feedback, correct = critic_modules[j]([taskInfo,
25             all_thinking[i][0], all_answer[i][0]],
26             critique_instruction)
27             all_feedback[i].append(feedback)
28
29 for i in range(len(FM_modules)):
30     refine_inputs = [taskInfo, all_thinking[i][0], all_answer[i]
31     ][0] + all_feedback[i]
32     thinking, answer = FM_modules[i](refine_inputs,
33     refine_instruction)
34     all_thinking[i].append(thinking)
35     all_answer[i].append(answer)
36
37 final_inputs = [taskInfo] + [all_thinking[i][1] for i in range(len(
38     FM_modules))] + [all_answer[i][1] for i in range(len(
39     FM_modules))]
40     thinking, answer = final_decision_module(final_inputs,
41     final_decision_instruction)
42
43 return answer

```

Code 5: Example discovered agent: Divide and Conquer Agent

```

1 def forward(self, taskInfo):
2     # Step 1: Decompose the problem into sub-problems
3     decomposition_instruction = "Please decompose the problem into
4     smaller, manageable sub-problems. List each sub-problem
5     clearly."
6     decomposition_module = FM_Module(['thinking', 'sub_problems'], '
7     Decomposition Module')
8
9     # Step 2: Assign each sub-problem to a specialized expert

```

```

7     sub_problem_instruction = "Please think step by step and then
      solve the sub-problem."
8     specialized_experts = [FM_Module(['thinking', 'sub_solution'], '
      Specialized Expert', role=role) for role in ['Physics Expert',
      'Chemistry Expert', 'Biology Expert', 'General Expert']]
9
10    # Step 3: Integrate the sub-problem solutions into the final
      answer
11    integration_instruction = "Given the solutions to the sub-problems
      , integrate them to provide a final answer to the original
      problem."
12    integration_module = FM_Module(['thinking', 'answer'], '
      Integration Module', temperature=0.1)
13
14    # Decompose the problem
15    thinking, sub_problems = decomposition_module([taskInfo],
      decomposition_instruction)
16
17    # Ensure sub_problems is a string and split into individual sub-
      problems
18    sub_problems_list = sub_problems.content.split('\n') if isinstance
      (sub_problems.content, str) else []
19
20    # Solve each sub-problem
21    sub_solutions = []
22    for i, sub_problem in enumerate(sub_problems_list):
23        sub_problem_info = Info('sub_problem', decomposition_module.
      __repr__(), sub_problem, i)
24        sub_thinking, sub_solution = specialized_experts[i % len(
      specialized_experts)]([sub_problem_info],
      sub_problem_instruction)
25        sub_solutions.append(sub_solution)
26
27    # Integrate the sub-problem solutions
28    integration_inputs = [taskInfo] + sub_solutions
29    thinking, answer = integration_module(integration_inputs,
      integration_instruction)
30
31    return answer

```

Code 6: Example discovered agent: Verified Multimodal Agent

```

1     def forward(self, taskInfo):
2         # Instruction for generating visual representation of the problem
3         visual_instruction = "Please create a visual representation (e.g.,
      diagram, graph) of the given problem."
4
5         # Instruction for verifying the visual representation
6         verification_instruction = "Please verify the accuracy and
      relevance of the visual representation. Provide feedback and
      suggestions for improvement if necessary."
7
8         # Instruction for solving the problem using the verified visual
      aid
9         cot_instruction = "Using the provided visual representation, think
      step by step and solve the problem."
10
11        # Instantiate the visual representation module, verification
      module, and Chain-of-Thought module
12        visual_module = FM_Module(['visual'], 'Visual Representation
      Module')
13        verification_module = FM_Module(['feedback', 'verified_visual'], '
      Verification Module')
14        cot_module = FM_Module(['thinking', 'answer'], 'Chain-of-Thought
      Module')

```

```
15
16     # Generate the visual representation of the problem
17     visual_output = visual_module([taskInfo], visual_instruction)
18     visual_representation = visual_output[0] # Using Info object
19         directly
20
21     # Verify the visual representation
22     feedback, verified_visual = verification_module([taskInfo,
23         visual_representation], verification_instruction)
24
25     # Use the verified visual representation to solve the problem
26     thinking, answer = cot_module([taskInfo, verified_visual],
27         cot_instruction)
28     return answer
```

K Cost of Experiments

A single run of search and evaluation on ARC (Section 3) costs approximately \$500 USD in OpenAI API costs, while a run within the reasoning and problem-solving domains (Section 3) costs about \$300 USD.

The primary expense comes from querying the “gpt-3.5-turbo-0125” model during the evaluation of discovered agents. Notably, the latest GPT-4 model, “gpt-4o-mini,” is less than one-third the price of “gpt-3.5-turbo-0125” and offers better performance, suggesting that we could achieve improved results with Meta Agent Search at just one-third of the cost. Additionally, as discussed in Section 4, the current naive evaluation function is both expensive and overlooks valuable information. We anticipate that future work adopting more sophisticated evaluation functions could significantly reduce the cost of ADAS algorithms.