

AUTOMATED DESIGN OF AGENTIC SYSTEMS

Anonymous authors

Paper under double-blind review

ABSTRACT

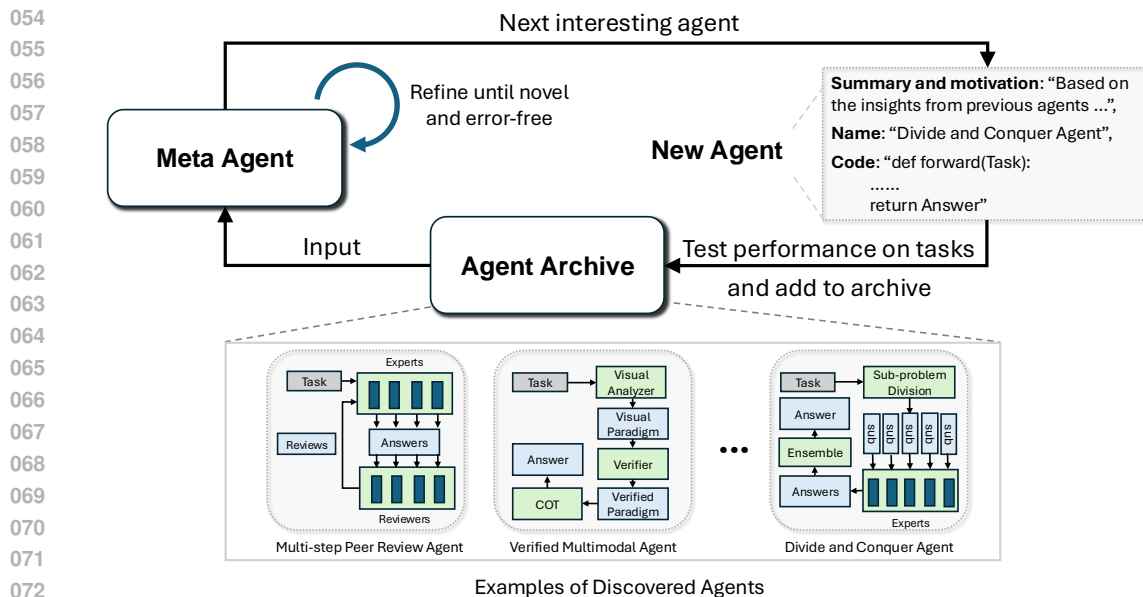
Researchers are investing substantial effort in developing powerful general-purpose agents, wherein Foundation Models are used as modules within *agentic systems* (e.g. Chain-of-Thought, Self-Reflection, Toolformer). However, the history of machine learning teaches us that hand-designed solutions are 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 most 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

1 INTRODUCTION

Foundation Models (FMs) such as GPT (OpenAI, 2024; 2022) and Claude (Anthropic, 2024b) are quickly being adopted as powerful general-purpose agents for agentic tasks that need flexible reasoning and planning (Wang et al., 2024). Despite recent advancements in FMs, solving problems reliably often requires an agent to be a compound agentic system with multiple components instead of a monolithic model query (Zaharia et al., 2024; Rocktäschel, 2024). Additionally, to enable agents to solve complex real-world tasks, they often need access to external tools such as search engines, code execution, and database queries. As a result, many effective building blocks of agentic systems have been proposed, such as chain-of-thought planning and reasoning (Wei et al., 2022; Yao et al., 2023; Hu & Clune, 2024), memory structures (Zhang et al., 2024c; Lewis et al., 2020), tool use (Schick et al., 2023; Qu et al., 2024), and self-reflection (Madaan et al., 2024; Shinn et al., 2023). Although these agents have already seen significant success across various applications (Wang et al., 2024), developing these building blocks and combining them into complex agentic systems often requires domain-specific manual tuning and substantial effort from both researchers and engineers.

However, the history of machine learning reveals a recurring theme: manually created artifacts become replaced by learned, more efficient solutions (Clune, 2019) over time as we get more compute and data (Sutton, 2019). An early example is from computer vision, where hand-designed features like HOG (Dalal & Triggs, 2005) were eventually replaced by learned features from Convolutional Neural Networks (CNNs, Krizhevsky et al. (2012)). More recently, AutoML methods (Hutter et al.,

¹All code will be open-sourced on release.



074
075
076
077
078

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. The detailed code of example agents can be found in Appendix H.

079
080
081
082
083
084
085
086
087
088
089

2019) and AI-Generating Algorithms (AI-GAs, Clune (2019)) have also demonstrated the superiority of learned AI systems compared to hand-designed AI systems. For example, the current best-performing CNN models come from Neural Architecture Search (Elsken et al., 2019; Shen et al., 2023) instead of manual design; in LLM alignment, learned loss functions (Lu et al., 2024a) outperform most hand-designed ones such as DPO (Rafailov et al., 2024); The AI Scientist (Lu et al., 2024b) demonstrates an automated research pipeline, including the development of novel ML algorithms; and an endless number of robotics learning environments can be automatically generated in works like OMNI-EPIC (Faldor et al., 2024), which demonstrate surprising creativity in generated environments and allow more efficient environment creation than the manual approach (see more examples in Section 5). Therefore, in this paper, we propose a new research question: *Can we automate the design of agentic systems?*

090
091
092
093
094
095
096
097
098
099
100

To explore the above research question, we describe a newly forming research area we call **Automated Design of Agentic Systems (ADAS)**, which aims to automatically invent novel building blocks and design powerful agentic systems (Section 2). We argue that ADAS may prove to be the fastest path to developing powerful agents, and show initial evidence that learned agents can greatly outperform hand-designed agents. Considering the tremendous number of building blocks yet to be discovered in agentic systems (Section 5), it would take a long time for our research community to discover them all. Even if we successfully discover most of the useful building blocks, combining them into effective agentic systems for massive real-world applications would still be challenging and time-consuming, given the many different ways the building blocks can combine and interact with each other. In contrast, with ADAS, the building blocks and agents can be learned in an automated fashion. ADAS may not only potentially save human effort in developing powerful agents but also could be a faster path to more effective solutions than manual design.

101
102
103
104
105
106
107

Although a few existing works can be considered as ADAS methods, most of them focus only on designing prompts (Yang et al., 2024; Fernando et al., 2024), greatly limiting their ability to invent flexible design patterns in agents (Section 5). In this paper, we show that there is an unexplored yet promising approach to ADAS where we can define the entire agentic system in code and new agents can be automatically discovered by a “meta” agent programming ever better ones in code. Given that most programming languages, such as Python, which we use in this paper, are Turing Complete (Boyer & Moore, 1983; Ladha, 2024), searching within a code space theoretically enables an ADAS algorithm to discover *any* possible agentic systems, including all components such as

prompts, tool use, workflows, and more. Furthermore, with recent FMs being increasingly proficient in coding, we can use FMs as meta agents to create new agents in code for ADAS, enabling novel agents to be programmed in an automated manner.

Following the aforementioned ideas, we present Meta Agent Search in this paper as one of the first algorithms in ADAS that enables complete design in code space (Figure 1). The core concept of Meta Agent Search is to instruct a meta agent to iteratively create interestingly new agents, evaluate them, add them to an archive that stores discovered agents, and use this archive to help the meta agent in subsequent iterations create yet more interestingly new agents. Similar to existing open-endedness algorithms that leverage human notions of interestingness (Zhang et al., 2024a; Lu et al., 2024c), we encourage the meta agent to explore interesting (e.g., novel or worthwhile) agents. To validate the proposed approach, we evaluate the proposed Meta Agent Search on: (1) the challenging ARC logic puzzle task (Chollet, 2019) that aims to test the general intelligence of an AI system, (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 4).

Our experiments show that the discovered agents substantially outperform state-of-the-art hand-designed baselines. For instance, our agents improve F1 scores on reading comprehension tasks in DROP (Dua et al., 2019) by **13.6/100** and accuracy rates on math tasks in MGSM (Shi et al., 2023) by **14.4%**. Additionally, they improve accuracy over baselines by **25.9%** and **13.2%** on GSM8K (Cobbe et al., 2021) and GSM-Hard (Gao et al., 2023) math tasks, respectively, *after transferring* across domains. The promising performance of our algorithm over hand-designed solutions illustrates the potential of ADAS in automating the design of agentic systems. Furthermore, the experiments demonstrate that the discovered agents not only perform well when transferring across similar domains but also exhibit strong performance when transferring across dissimilar domains, such as from mathematics to reading comprehension. This highlights the robustness and transferability of the agentic systems discovered by Meta Agent Search. In conclusion, our work opens up many exciting research directions and encourages further studies (Section 6).

2 AUTOMATED DESIGN OF AGENTIC SYSTEMS (ADAS)

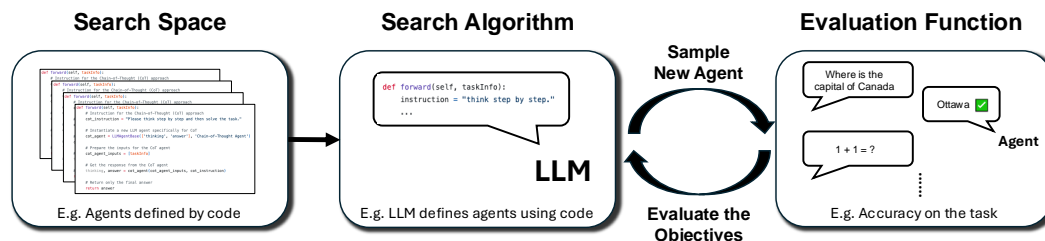


Figure 2: **The three key components of Automated Design of Agentic Systems (ADAS).** The search space determines which agentic systems can be represented in ADAS. The search algorithm specifies how the ADAS method explores the search space. The evaluation function defines how to evaluate a candidate 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 (Chase, 2024; Ng, 2024). In this paper, we describe a newly forming research area Automated Design of Agentic Systems (ADAS). Similar to research areas in AI-GAs (Clune, 2019) and AutoML (Hutter et al., 2019), such as Neural Architecture Search (Elsken et al., 2019), we formulate ADAS as an optimization process and identify three key components of ADAS algorithms (Figure 2).

Formulation

Automated Design of Agentic Systems (ADAS) involves using a **search algorithm** to discover agentic systems across a **search space** that **optimize an evaluation function**.

- **Search Space:** The search space defines which agentic systems can be represented and thus discovered in ADAS. For example, works like PromptBreeder (Fernando et al., 2024) mutate only the text prompts of an agent, but their other components, such as workflow, remain the same. Thus, in these search spaces, agents that have a different workflow than the predefined one can not be represented. Existing works also explore search spaces such as graph structures (Zhuge et al., 2024) and feed-forward networks (Liu et al., 2023).
- **Search Algorithm:** The search algorithm defines how ADAS algorithms explore the search space. Since the search space is often very large or even unbounded, the exploration-exploitation trade-off (Sutton & Barto, 2018) should be considered. Ideally, the algorithm can both quickly discover high-performance agentic systems and avoid remaining stuck in a local optimum. Existing approaches include using Reinforcement Learning (Zhuge et al., 2024) or an FM iteratively generating new solutions (Fernando et al., 2024) as search algorithms.
- **Evaluation Function:** Depending on the application of the ADAS algorithm, we may consider different objectives to optimize, such as performance, cost, latency, or safety of agents. An evaluation function defines how to evaluate a candidate agent on those objectives. For example, to assess the agent’s performance on unseen future data, a simple method is to calculate the accuracy rate on the validation data for a task, which is commonly adopted in existing works (Zhuge et al., 2024; Fernando et al., 2024).

Although many search space designs are possible and some have already been explored (Section 5), 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 ever 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 (Liu et al., 2023) or graphs (Zhuge et al., 2024), 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 (LangChainAI, 2022) 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.

3 OUR ALGORITHM: META AGENT SEARCH

In this section, we present Meta Agent Search, a simple yet effective algorithm to demonstrate the approach of defining and searching for agents in code. The core idea of Meta Agent Search is to adopt FMs as meta agents to iteratively program interestingly new agents based on an ever-growing archive of previous discoveries. Although any possible building blocks and agentic systems can theoretically be programmed by the meta agent from scratch, it is inefficient in practice to avoid providing the meta agent any basic functions such as FM query APIs or existing tools. Therefore, in this paper, we define a simple framework (within 100 lines of code) for the meta agent, providing it with a basic set of essential functions like querying FMs or formatting prompts. As a result, the meta agent only needs to program a “forward” function to define a new agentic system, similar to the practice in FunSearch (Romera-Paredes et al., 2024). This function takes in the information of the task and outputs the agent’s response to the task. Details of the framework codes and examples of the agents defined with this framework can be found in Appendix D.

As shown in Figure 1, the core idea of Meta Agent Search is to have a meta agent iteratively program new agents in code. The algorithm proceeds as follows: (1) The archive is (optionally) initialized with baseline agents such as Chain-of-Thought (Wei et al., 2022) and Self-Refine (Madaan et al., 2024; Shinn et al., 2023). (2) Conditioned on the archive, the meta agent designs a new agent by generating a high-level description of the new idea for an agent system and then implementing it in code. The design then undergoes two self-reflection (Madaan et al., 2024; Shinn et al., 2023) steps by the meta agent to ensure it is novel. (3) The generated agent is evaluated using validation data from the target domain. If errors occur during evaluation, the meta agent performs a self-reflection step to

refine the design, repeating this process up to five times if necessary. (4) Finally, the agent is added to the archive along with its evaluation metrics, and the process continues with the updated archive until the maximum number of iterations is reached. A pseudocode of the algorithm is provided in Appendix I.

Similar to existing open-endedness algorithms that leverage human notions of interestingness (Zhang et al., 2024a; Lu et al., 2024c), we encourage the meta agent to explore interestingly new (e.g., novel or worthwhile) agents based on an ever-growing archive of previous discoveries. Here, we calculate the performance (e.g., success rate or F1 score) as the metrics for the meta agent to maximize. The prompt and more details are presented in Appendix C.

4 EXPERIMENTS

We conduct extensive experiments on: (1) the ARC challenge (Chollet, 2019) (Section 4.1), (2) four popular benchmarks assessing the agent’s abilities on reading comprehension, math, science questions, and multi-task problem solving (Section 4.2), and (3) the transferability of discovered agents on math to held-out math tasks and non-math tasks (Section 4.3). We use an identical implementation of the algorithm across different tasks, with the only variation being task-specific descriptive text included in the prompt (details are available in Appendix C). Across all experiments, we find that the discovered agents substantially outperform baseline state-of-the-art hand-designed agents and maintain superior performance even when transferred across domains and models.

4.1 CASE STUDY: ARC CHALLENGE

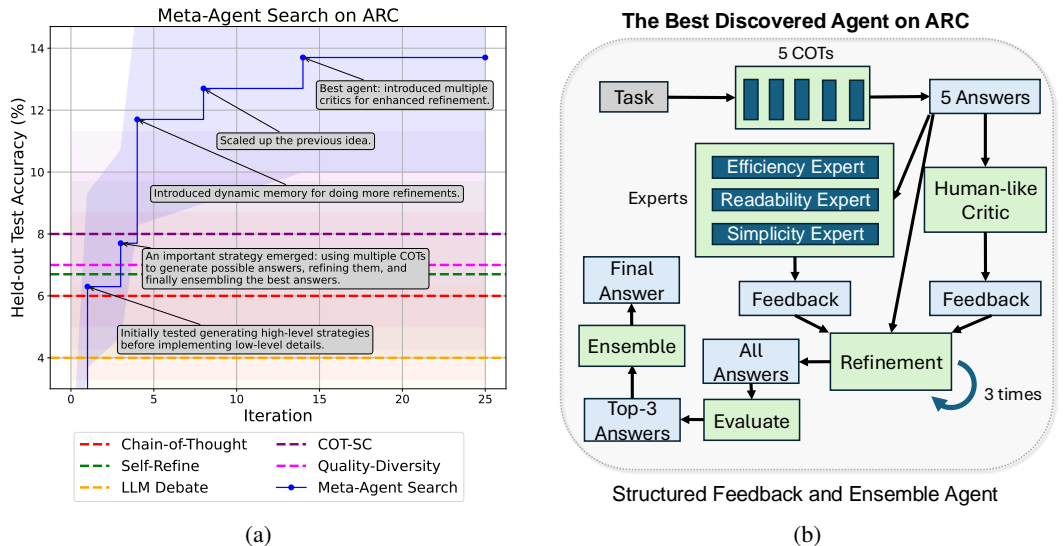


Figure 3: **The results of Meta Agent Search on the ARC challenge.** (a) Meta Agent Search progressively discovers high-performance agents based on an ever-growing archive of previous discoveries. We report the median accuracy and the 95% bootstrap confidence interval on a held-out test set by evaluating agents five times. (b) The visualization of the best agent discovered by Meta Agent Search on the ARC challenge. Detailed implementation of this agent is available in Appendix E.

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 (Chollet, 2019). This 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.

Setup. Following common practice (Greenblatt, 2024), we require the agent to write code for the transformation rule instead of answering directly. We provide tool functions in the framework (described in Section 3) 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 3. Meta Agent Search runs for 25 iterations and the meta agent uses GPT-4 (OpenAI, 2024), while discovered agents and baselines are evaluated using GPT-3.5 (OpenAI, 2022) to reduce compute cost. More algorithmic details and examples of ARC questions can be found in Appendix E.

Baselines. We compared against five state-of-the-art hand-designed agents: (1) Chain-of-Thought (COT, Wei et al. (2022)), 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, Wang et al. (2023b)), which ensembles multiple parallel answers from COT to produce a more accurate answer; (3) Self-Refine (Madaan et al., 2024; Shinn et al., 2023), which allows iterative self-reflection to correct mistakes made in previous attempts; (4) LLM-Debate (Du et al., 2023), 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 (Lu et al., 2024c), which produces and ensembles diverse answers to better explore potential solutions. The selected baselines represent widely adopted agent designs in the agent literature, embodying key design patterns and approaches frequently utilized across various applications. By “state-of-the-art,” we refer to these baseline designs as exemplifying important advancements and practices within the field. We also use all baselines as initial seeds in the archive for Meta Agent Search, with additional results for empty initialization provided in Appendix J. To ensure fair comparisons, all baseline implementations were developed using the same framework as the Meta Agent, providing a consistent and equitable evaluation environment. More details about baselines can be found in Appendix G.

Results and Analysis. As shown in Figure 3a, Meta Agent Search effectively and progressively discovers agents that perform better than state-of-the-art hand-designed baselines. Important breakthroughs are highlighted in the text boxes. As is critical in prior works on open-endedness and AI-GAs (Zhang et al., 2024a; Faldor et al., 2024; Wang et al., 2019; 2020; Lehman & Stanley, 2011), Meta Agent Search innovates based on a growing archive of previous stepping stones. For example, an important design pattern emerged in iteration 3 where it uses multiple COTs to generate possible answers, refines them, and finally ensembles the best answers. This became a crucial stepping stone that subsequent designs tended to utilize. Additionally, the best-discovered agent is shown in Figure 3b, where a complex feedback mechanism is adopted to refine answers more effectively. Careful observation of the search progress reveals that this sophisticated feedback mechanism did not appear suddenly. Instead, the ideas of incorporating diverse feedback, evaluating for various specific traits (via experts) such as efficiency and simplicity, and simulating human-like feedback emerged in iterations 5, 11, and 12, respectively. The final mechanism is an innovation based on these three stepping stones. This illustrates that even though these stepping stones did not achieve high performance immediately upon emergence, later discoveries benefited from these innovations by combining different stepping stones, resembling crossover in evolution via LLMs (Meyerson et al., 2023). Overall, the results showcase the potential of ADAS and the effectiveness of Meta Agent Search to progressively discover agents that outperform state-of-the-art hand-designed baselines and invent novel design patterns through the innovation and combination of stepping stones.

4.2 REASONING AND PROBLEM-SOLVING DOMAINS

Setup. Next, we investigate the potential of our algorithm to improve the capabilities of agents across math, reading, and reasoning domains. We test Meta Agent Search on four popular benchmarks: (1) DROP (Dua et al., 2019) for evaluating **Reading Comprehension**; (2) MGSM (Shi et al., 2023) for evaluating **Math** capability under a multi-lingual setting; (3) MMLU (Hendrycks et al., 2021) for evaluating **Multi-task** Problem Solving; and (4) GPQA (Rein et al., 2023) 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 (OpenAI, 2024), while the discovered agents and baselines are evaluated using GPT-3.5 (OpenAI, 2022). More details about datasets and experiment settings can be found in Appendix F.

Baselines. We adopt all baselines introduced in Section 4.1. Additionally, since the above domains require strong reasoning skills, we include two additional baselines that specifically focus on enhancing the reasoning capabilities of agents for a more thorough comparison: (1) Step-back Abstraction (Zheng et al., 2023), which instructs agents to first consider the principles involved in solving the task for better reasoning; (2) Role Assignment (Xu et al., 2023), which assigns different roles to FMs to obtain better answers. Furthermore, we compare our approach with the state-of-the-art prompt optimization baseline OPRO (Yang et al., 2024) to highlight the advantages of learning all possible components of agents rather than focusing solely on prompts. More details about the baselines can be found in Appendix G.

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. We report the test accuracy and the 95% bootstrap confidence interval on held-out test sets. The search is conducted independently for each domain. Here, and in all tables below, we bold the entry with the highest performance for each domain, as well as all entries whose median falls within the 95% confidence interval of the highest-performing treatment.

Agent Name	F1 Score		Accuracy (%)		
	Reading Comprehension	Math	Multi-task	Science	
State-of-the-art Hand-designed Agents					
Chain-of-Thought (Wei et al., 2022)	64.2 ± 0.9	28.0 ± 3.1	65.4 ± 3.3	29.2 ± 3.1	
COT-SC (Wang et al., 2023b)	64.4 ± 0.8	28.2 ± 3.1	65.9 ± 3.2	30.5 ± 3.2	
Self-Refine (Madaan et al., 2024)	59.2 ± 0.9	27.5 ± 3.1	63.5 ± 3.4	31.6 ± 3.2	
LLM Debate (Du et al., 2023)	60.6 ± 0.9	39.0 ± 3.4	65.6 ± 3.3	31.4 ± 3.2	
Step-back Abstraction (Zheng et al., 2023)	60.4 ± 1.0	31.1 ± 3.2	65.1 ± 3.3	26.9 ± 3.0	
Quality-Diversity (Lu et al., 2024c)	61.8 ± 0.9	23.8 ± 3.0	65.1 ± 3.3	30.2 ± 3.1	
Role Assignment (Xu et al., 2023)	65.8 ± 0.9	30.1 ± 3.2	64.5 ± 3.3	31.1 ± 3.1	
Automated Design of Agentic Systems on Different Domains					
Prompt Optimization (Yang et al., 2024)	69.1 ± 0.9	30.6 ± 3.2	67.6 ± 3.2	32.9 ± 3.2	
Meta Agent Search (Ours)	79.4 ± 0.8	53.4 ± 3.5	69.6 ± 3.2	34.6 ± 3.2	

Results and Analysis. The results across multiple domains demonstrate that Meta Agent Search can discover agents that outperform state-of-the-art hand-designed agents (Table 1). We want to highlight the substantial gap between the learned agents and hand-designed agents in the Reading Comprehension and Math domains, with improvements in F1 scores by **13.6/100** and accuracy rates by **14.4%**, respectively. While Meta Agent Search also outperforms baselines in the Multi-task and Science domains, the gap is smaller. We hypothesize that for challenging questions in the Science and Multi-task domains, the knowledge in FMs is not sufficient to solve the questions, limiting the improvement through optimizing agentic systems, which is a problem that will diminish as FMs improve. In contrast, in the Reading Comprehension and Math domains, FMs possess adequate knowledge to solve the questions, and errors could mainly be hallucinations or calculation mistakes, which can be mitigated through well-designed agentic systems, like the ones discovered by Meta Agent Search. Additionally, when compared to prompt optimization methods, the results demonstrate that our proposed Meta Agent Search consistently outperforms them across all domains. This comparison further strengthens our argument that defining agents in code and enabling the learning of all components offer significant advantages. Overall, the results across various domains showcase the effectiveness of Meta Agent Search in searching for agents tailored to specific domains. This could be increasingly useful for saving human efforts and developing better task-specific agents as we continue to create agents for a diverse set of applications (Wang et al., 2024).

4.3 GENERALIZATION AND TRANSFERABILITY

In the previous sections, we illustrated that Meta Agent Search can find effective agents for individual tasks. In this section, we further demonstrate the transferability and generalizability of the discovered agents. To demonstrate the generalizability of the invented building blocks and de-

sign patterns, we transfer discovered agents from the MGSM (Math) domain to both math and non-math domains to test their ability to generalize across different tasks. We evaluate the top 3 agents from MGSM by transferring them to (1) popular math domains: GSM8K (Cobbe et al., 2021), GSM-Hard (Gao et al., 2023), and (2) non-math domains: MMLU (Multi-task) and DROP (Reading Comprehension), as detailed in Section 4.2. As shown in Table 2, Meta Agent Search consistently outperforms the baselines. Notably, our agents improve accuracy by **25.9%** on GSM8K and **13.2%** on GSM-Hard compared to the baselines when transferring within math domains. More surprisingly, we find that agents discovered in the math domain can also be transferred to non-math domains. While their performance does not fully match agents specifically designed for the target domains, they still outperform state-of-the-art hand-designed baselines. More results of transfers across domains are shown in Appendix B.

We also observe similar superiority when transferring agents across different FMs on ARC. We test the top 3 agents with the best test accuracy evaluated with GPT-3.5 on ARC and then transfer them to Claude-Haiku (Anthropic, 2024a), GPT-4 (OpenAI, 2024), and Claude-Sonnet (Anthropic, 2024b). 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. These results on transferring across domains and models highlight Meta Agent Search’s ability to discover generalizable design patterns and agentic systems.

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. We report the test accuracy and the 95% bootstrap confidence interval. The names of the top agents are generated by Meta Agent Search.

Agent Name	Accuracy (%)				F1 Score	
	MGSM	GSM8K	GSM-Hard	MMLU	DROP	
Manually Designed Agents						
Chain-of-Thought (Wei et al., 2022)	28.0 ± 3.1	34.9 ± 3.2	15.0 ± 2.5	65.4 ± 3.3	64.2 ± 0.9	
COT-SC (Wang et al., 2023b)	28.2 ± 3.1	37.8 ± 3.4	15.5 ± 2.5	65.9 ± 3.2	64.4 ± 0.8	
Self-Refine (Madaan et al., 2024)	27.5 ± 3.1	38.9 ± 3.4	15.1 ± 2.4	63.5 ± 3.4	59.2 ± 0.9	
LLM Debate (Du et al., 2023)	39.0 ± 3.4	43.6 ± 3.4	17.4 ± 2.6	65.6 ± 3.3	60.6 ± 0.9	
Step-back Abstraction (Zheng et al., 2023)	31.1 ± 3.2	31.5 ± 3.3	12.2 ± 2.3	65.1 ± 3.3	60.4 ± 1.0	
Quality-Diversity (Lu et al., 2024c)	23.8 ± 3.0	28.0 ± 3.1	14.1 ± 2.4	65.1 ± 3.1	61.8 ± 0.9	
Role Assignment (Xu et al., 2023)	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	

5 RELATED WORK

Agentic Systems. Researchers develop various building blocks and design patterns for different applications. Important building blocks for agentic systems include: prompting techniques (Chen et al., 2023a; Schulhoff et al., 2024), chain-of-thought-based planning and reasoning methods (Wei et al., 2022; Yao et al., 2023; Hu & Clune, 2024), reflection (Madaan et al., 2024; Shinn et al., 2023), developing new skills for embodied agents in code (Wang et al., 2023a; Vemprala et al., 2023), external memory and RAG (Zhang et al., 2024c; Lewis et al., 2020), tool use (Qu et al., 2024; Schick et al., 2023; Nakano et al., 2021), assigning FM modules in the agentic system with different roles and enabling them to collaborate (Hong et al., 2023; Wu et al., 2023; Qian et al., 2023; Xu et al., 2023; Qian et al., 2024), and enabling the agent to instruct itself for the next action (Richards, 2023), 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

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 (Wei et al., 2022)	6.0 ± 2.7	4.3 ± 2.2	17.7 ± 4.4	25.3 ± 5.0
COT-SC (Wang et al., 2023b)	8.0 ± 3.2	5.3 ± 2.5	19.7 ± 4.5	26.3 ± 4.9
LLM Debate (Du et al., 2023)	4.0 ± 2.2	1.7 ± 1.5	19.0 ± 4.5	24.7 ± 4.8
Self-Refine (Madaan et al., 2024)	6.7 ± 2.7	6.3 ± 2.8	23.0 ± 5.2	39.3 ± 5.5
Quality-Diversity (Lu et al., 2024c)	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

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.

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 (Yang et al., 2024; Fernando et al., 2024; Zhou et al., 2024a). However, these prompts are often domain-specific and difficult to generalize. Some works optimize role definitions within prompts (Yuan et al., 2024; Chen et al., 2023c;b; Wu et al., 2023), as assigning personas or roles to agents has been shown to be beneficial (Xu et al., 2023). Although tuning prompts can improve performance, 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. In these formulations, the FM with a certain prompt is considered a transformation function for text on nodes, and the information flow of the text is considered as edges. For example, DyLAN (Liu et al., 2023) uses FMs to optimize connections between nodes in a network, DSPy (Khattab et al., 2024) optimizes across the Cartesian product of a set of possible nodes, and GPT-Swarm (Zhuge et al., 2024) uses reinforcement learning to optimize node connections. Although these approaches optimize workflows, many components like tool usage remain fixed. AgentOptimizer (Zhang et al., 2024b) learns the tools used in agents and Agent Symbolic Learning (Zhou et al., 2024b) attempts to learn prompts, tools, and workflows together. While Agent Symbolic Learning shares similar motivations to learn more components in agents, it manually designs the search space for each component separately, which may make it a harder search space for search algorithms. In addition, it improves agents based on an existing complex agent, without showing the emergence of new design patterns or building blocks. In contrast, our work represents all components in code, allowing the search to be easier by leveraging human efforts in the existing codebase of agents and FMs’ expertise in coding. We also demonstrate how novel and diverse building blocks and design patterns emerge from a set of basic agent designs, illustrating the potential creativity that can emerge from ADAS.

We also include additional related work in Appendix A.1.

6 DISCUSSION AND CONCLUSION

Safety Considerations. 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 (Rokon et al., 2020; Chen et al., 2021). 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

486 includes the proposed Automated Design of Agentic Systems (ADAS) as a new area in AI-GA
487 research, which could potentially contribute to an even faster way to create Artificial General Intelli-
488 gence (AGI) than the current manual approach (Clune, 2019). The question of whether and why we
489 should pursue AGI and AI-GA has been discussed in many papers (Clune, 2019; Ecoffet et al., 2020;
490 Bostrom, 2002; Yudkowsky et al., 2008; Bengio et al., 2024), and is beyond the scope of this paper.
491 Specifically as regards to ADAS, we believe it is net beneficial to publish this work. First, this work
492 demonstrates that with the available API access to powerful FMs, it is easy to program powerful
493 ADAS algorithms, and do so without any expensive hardware like GPUs. We feel it is beneficial to
494 let the community know such algorithms are powerful and easy to create, so they can be informed
495 and account for them. Moreover, by sharing this information, we hope to motivate follow-up work
496 into safe-ADAS, such as algorithms that conduct ADAS safely during both search itself (e.g. not
497 risking running any harmful code) and that refuse to create dishonest, unhelpful, and/or harmful
498 agents. Such an open-source research approach to create safe-ADAS could be a better way to create
499 safer AI systems (Caldwell, 2011; Meta, 2024). One direction we find particularly promising is
500 to simply ask the Meta Agent Search algorithm to be safe during training and only create helpful,
501 harmless, honest agents, potentially incorporating ideas such as Constitutional AI (Bai et al., 2022).

502 **Future Work.** Our work also opens up many future research directions. Below, we discuss a few,
503 with additional directions provided in Appendix A.2.

- 504 • **Higher-order ADAS.** Since the meta agent used in ADAS to program new agents in code is
505 also an agent, ADAS can become self-referential where the meta agent can be improved through
506 ADAS as well. It would be an exciting direction to have a higher order of meta-learning to allow
507 the learning of the meta agent and even the meta-meta agent, etc. (Lu et al., 2023; Schmidhuber,
508 1987)
- 509 • **Online Continual Learning.** As agents are deployed, they will receive vast amounts of feedback
510 from both task environments and users. Continuously improving agents based on this extensive
511 feedback is challenging for human developers. However, with ADAS automating the design and
512 enhancement of agents, online continual learning becomes feasible post-deployment.
- 513 • **Multi-objective ADAS.** We only consider one objective (i.e., performance) to optimize in this
514 paper, but in practice, multiple objectives are often considered, such as cost, latency, and robust-
515 ness of agentic systems (Hu et al., 2021; Huang et al., 2023). Thus, integrating multi-objective
516 search algorithms (Deb et al., 2002) in ADAS could be promising.
- 517 • **Towards a Better Understanding of FMs.** Works from Neural Architecture Search (Huang
518 et al., 2023) show that by observing the emerged architecture, we could gain more insights into
519 Neural Networks. In this paper, we also gained insights about FMs from the results. For exam-
520 ple, the best agent with GPT-3.5 involves a complex feedback mechanism, but when we transfer
521 to other advanced models, the agent with a simpler feedback mechanism but more refinement
522 becomes a better agent (Section 4.3). This shows that GPT-3.5 may have a worse capability in
523 evaluating and refining the answers, so it needs a complex feedback mechanism for better refine-
524 ment, while other advanced models benefit more from a simpler feedback mechanism.

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

REFERENCES

- 540
541
542 Anthropic. Introducing the next generation of claude. <https://www.anthropic.com/news/claude-3-family>, March 2024a. Blog post.
- 543
544 Anthropic. Introducing claude 3.5 sonnet. <https://www.anthropic.com/news/claude-3-5-sonnet>, June 2024b. Blog post.
- 545
546
547 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,
548 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harm-
549 lessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- 550
551 Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Trevor Darrell, Yu-
552 val Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, et al. Managing extreme ai risks
553 amid rapid progress. *Science*, 384(6698):842–845, 2024.
- 554
555 N Bostrom. Existential Risks: analyzing human extinction scenarios and related hazards. *Journal*
556 *of Evolution and Technology*, 9, 2002.
- 557
558 Robert S Boyer and J Strother Moore. *A mechanical proof of the Turing completeness of pure LISP*.
Citeseer, 1983.
- 559
560 Tracey Caldwell. Ethical hackers: putting on the white hat. *Network Security*, 2011(7):10–13, 2011.
- 561
562 Harrison Chase. What is an agent? [https://blog.langchain.dev/](https://blog.langchain.dev/what-is-an-agent/)
563 [what-is-an-agent/](https://blog.langchain.dev/what-is-an-agent/), June 2024. Blog post.
- 564
565 Banghao Chen, Zhaofeng Zhang, Nicolas Langrené, and Shengxin Zhu. Unleashing the poten-
566 tial of prompt engineering in large language models: a comprehensive review. *arXiv preprint*
arXiv:2310.14735, 2023a.
- 567
568 Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Sesay Jaward, Karlsson Börje, Jie Fu, and Yemin
569 Shi. Autoagents: The automatic agents generation framework. *arXiv preprint*, 2023b.
- 570
571 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared
572 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 573
574 Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu,
575 Yaxi Lu, Yi-Hsin Hung, Chen Qian, et al. Agentverse: Facilitating multi-agent collaboration and
576 exploring emergent behaviors. In *The Twelfth International Conference on Learning Representa-*
577 *tions*, 2023c.
- 578
579 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li,
580 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.
- 581
582 François Chollet. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.
- 583
584 Jeff Clune. Ai-gas: Ai-generating algorithms, an alternate paradigm for producing general artificial
585 intelligence. *arXiv preprint arXiv:1905.10985*, 2019.
- 586
587 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
588 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to
solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 589
590 Antoine Cully and Yiannis Demiris. Quality and diversity optimization: A unifying modular frame-
591 work. *IEEE Transactions on Evolutionary Computation*, 22(2):245–259, 2017.
- 592
593 N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE Com-*
puter Society Conference on Computer Vision and Pattern Recognition (CVPR'05), volume 1, pp.
886–893 vol. 1, 2005. doi: 10.1109/CVPR.2005.177.

- 594 Kalyanmoy Deb, Amrit Pratap, Sameer Agarwal, and TAMT Meyarivan. A fast and elitist mul-
595 tiobjective genetic algorithm: Nsga-ii. *IEEE transactions on evolutionary computation*, 6(2):
596 182–197, 2002.
- 597
- 598 Aaron Dharna, Amy K Hoover, Julian Togelius, and Lisa B Soros. Transfer dynamics in emergent
599 evolutionary curricula. *IEEE Transactions on Games*, 15(2):157–170, 2022.
- 600
- 601 Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improv-
602 ing factuality and reasoning in language models through multiagent debate. *arXiv preprint*
603 *arXiv:2305.14325*, 2023.
- 604
- 605 Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner.
606 DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In
607 Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of*
608 *the North American Chapter of the Association for Computational Linguistics: Human Language*
609 *Technologies, Volume 1 (Long and Short Papers)*, pp. 2368–2378, Minneapolis, Minnesota, June
2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1246.
- 610
- 611 Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel. RL²: Fast
612 reinforcement learning via slow reinforcement learning. In *International Conference on Learning*
613 *Representations*, 2017.
- 614
- 615 Adrien Ecoffet, Jeff Clune, and Joel Lehman. Open questions in creating safe open-ended AI:
616 Tensions between control and creativity. In *Conference on Artificial Life*, pp. 27–35. MIT Press,
2020.
- 617
- 618 Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey.
619 *Journal of Machine Learning Research*, 20(55):1–21, 2019.
- 620
- 621 Maxence Faldor, Jenny Zhang, Antoine Cully, and Jeff Clune. Omni-epic: Open-endedness via mod-
622 els of human notions of interestingness with environments programmed in code. *arXiv preprint*
623 *arXiv:2405.15568*, 2024.
- 624
- 625 Chrisantha Fernando, Dylan Sunil Banarse, Henryk Michalewski, Simon Osindero, and Tim
Rocktäschel. Promptbreeder: Self-referential self-improvement via prompt evolution, 2024.
- 626
- 627 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation
628 of deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- 629
- 630 Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and
631 Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine*
Learning, pp. 10764–10799. PMLR, 2023.
- 632
- 633 Ryan Greenblatt. Getting 50% sota on arc-agi with gpt-4. [https://redwoodresearch.](https://redwoodresearch.substack.com/p/getting-50-sota-on-arc-agi-with-gpt)
634 [substack.com/p/getting-50-sota-on-arc-agi-with-gpt](https://redwoodresearch.substack.com/p/getting-50-sota-on-arc-agi-with-gpt), July 2024. Techni-
635 cal Report.
- 636
- 637 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
638 Steinhardt. Measuring massive multitask language understanding. In *International Conference*
on Learning Representations, 2021.
- 639
- 640 Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang,
641 Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-
642 agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 2023.
- 643
- 644 Shengran Hu and Jeff Clune. Thought Cloning: Learning to think while acting by imitating human
645 thinking. *Advances in Neural Information Processing Systems*, 36, 2024.
- 646
- 647 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*,
pp. 1–10, 2021.

- 648 Shihua Huang, Zhichao Lu, Kalyanmoy Deb, and Vishnu Naresh Boddeti. Revisiting residual net-
649 works for adversarial robustness. In *Proceedings of the IEEE/CVF Conference on Computer*
650 *Vision and Pattern Recognition*, pp. 8202–8211, 2023.
- 651 Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren. *Automated machine learning: methods, sys-*
652 *tems, challenges*. Springer Nature, 2019.
- 654 Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Saiful
655 Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, Heather Miller, et al. Dspy: Compil-
656 ing declarative language model calls into state-of-the-art pipelines. In *The Twelfth International*
657 *Conference on Learning Representations*, 2024.
- 658 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-
659 lutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- 661 Abraham Ladha. Lecture 11: Turing-completeness. [https://faculty.cc.gatech.edu/](https://faculty.cc.gatech.edu/~ladha/S24/4510/L11.pdf)
662 [~ladha/S24/4510/L11.pdf](https://faculty.cc.gatech.edu/~ladha/S24/4510/L11.pdf), 2024. CS 4510 Automata and Complexity, February 21st,
663 2024, Scribed by Rishabh Singhal.
- 664 LangChainAI. Langchain: Build context-aware reasoning applications. [https://github.](https://github.com/langchain-ai/langchain)
665 [com/langchain-ai/langchain](https://github.com/langchain-ai/langchain), 2022.
- 666 Joel Lehman and Kenneth O Stanley. Abandoning objectives: Evolution through the search for
667 novelty alone. *Evolutionary computation*, 19(2):189–223, 2011.
- 668 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
669 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented genera-
670 tion for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:
671 9459–9474, 2020.
- 672 Fei Liu, Tong Xialiang, Mingxuan Yuan, Xi Lin, Fu Luo, Zhenkun Wang, Zhichao Lu, and Qingfu
673 Zhang. Evolution of heuristics: Towards efficient automatic algorithm design using large language
674 model. In *Forty-first International Conference on Machine Learning*, 2024.
- 675 Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. Dynamic llm-agent network: An llm-
676 agent collaboration framework with agent team optimization. *arXiv preprint arXiv:2310.02170*,
677 2023.
- 678 Chris Lu, Sebastian Towers, and Jakob Foerster. Arbitrary order meta-learning with simple
679 population-based evolution. In *ALIFE 2023: Ghost in the Machine: Proceedings of the 2023*
680 *Artificial Life Conference*. MIT Press, 2023.
- 681 Chris Lu, Samuel Holt, Claudio Fanconi, Alex J Chan, Jakob Foerster, Mihaela van der Schaar, and
682 Robert Tjarko Lange. Discovering preference optimization algorithms with and for large language
683 models. *arXiv preprint arXiv:2406.08414*, 2024a.
- 684 Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The AI Scien-
685 tist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*,
686 2024b.
- 687 Cong Lu, Shengran Hu, and Jeff Clune. Intelligent go-explore: Standing on the shoulders of giant
688 foundation models. *arXiv preprint arXiv:2405.15143*, 2024c.
- 689 Zhichao Lu, Ian Whalen, Vishnu Boddeti, Yashesh Dhebar, Kalyanmoy Deb, Erik Goodman, and
690 Wolfgang Banzhaf. Nsga-net: neural architecture search using multi-objective genetic algorithm.
691 In *Proceedings of the genetic and evolutionary computation conference*, pp. 419–427, 2019.
- 692 Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayara-
693 man, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via
694 coding large language models. In *The Twelfth International Conference on Learning Representa-*
695 *tions*, 2023.

- 702 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri
703 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
704 with self-feedback. *Advances in Neural Information Processing Systems*, 36, 2024.
- 705
706 Meta. Open source ai is the path forward. [https://about.fb.com/news/2024/07/
707 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.
- 708 Elliot Meyerson, Mark J Nelson, Herbie Bradley, Adam Gaier, Arash Moradi, Amy K Hoover, and
709 Joel Lehman. Language model crossover: Variation through few-shot prompting. *arXiv preprint
710 arXiv:2302.12170*, 2023.
- 711 Shen-yun Miao, Chao-Chun Liang, and Keh-Yih Su. A diverse corpus for evaluating and developing
712 english math word problem solvers. In *Proceedings of the 58th Annual Meeting of the Association
713 for Computational Linguistics*, pp. 975–984, 2020.
- 714
715 Jean-Baptiste Mouret and Jeff Clune. Illuminating search spaces by mapping elites. *arXiv preprint
716 arXiv:1504.04909*, 2015.
- 717 Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christo-
718 pher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted
719 question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- 720
721 Andrew Ng. Issue 253. <https://www.deeplearning.ai/the-batch/issue-253/>,
722 June 2024. Newsletter issue.
- 723 Ben Norman and Jeff Clune. First-explore, then exploit: Meta-learning intelligent exploration. *arXiv
724 preprint arXiv:2307.02276*, 2023.
- 725
726 OpenAI. Introducing chatgpt. <https://openai.com/index/chatgpt/>, November 2022.
727 Blog post.
- 728 OpenAI. Simple evals, 2023. URL <https://github.com/openai/simple-evals>. Ac-
729 cessed: 2024-08-10.
- 730
731 OpenAI. Gpt-4 technical report, 2024.
- 732 Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and
733 Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings
734 of the 36th annual acm symposium on user interface software and technology*, pp. 1–22, 2023.
- 735
736 Arkil Patel, Satwik Bhattamishra, and Navin Goyal. Are NLP models really able to solve simple
737 math word problems? In *Proceedings of the 2021 Conference of the North American Chapter
738 of the Association for Computational Linguistics: Human Language Technologies*, pp. 2080–
739 2094, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.
naacl-main.168.
- 740
741 Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu,
742 and Maosong Sun. Communicative agents for software development. *arXiv preprint
743 arXiv:2307.07924*, 2023.
- 744
745 Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang,
746 Zhiyuan Liu, and Maosong Sun. Scaling large-language-model-based multi-agent collaboration.
arXiv preprint arXiv:2406.07155, 2024.
- 747
748 Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu,
749 and Ji-Rong Wen. Tool learning with large language models: A survey. *arXiv preprint
arXiv:2405.17935*, 2024.
- 750
751 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
752 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances
753 in Neural Information Processing Systems*, 36, 2024.
- 754
755 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.

- 756 Toran Bruce Richards. Autogpt. [https://github.com/Significant-Gravitas/](https://github.com/Significant-Gravitas/AutoGPT)
757 `AutoGPT`, 2023. GitHub repository.
- 758
- 759 Tim Rocktäschel. *Artificial Intelligence: 10 Things You Should Know*. Seven Dials, September
760 2024. ISBN 978-1399626521.
- 761
- 762 Md Omar Faruk Rokon, Risul Islam, Ahmad Darki, Evangelos E Papalexakis, and Michalis Falout-
763 sos. SourceFinder: Finding malware Source-Code from publicly available repositories in GitHub.
764 In *23rd International Symposium on Research in Attacks, Intrusions and Defenses (RAID 2020)*,
765 pp. 149–163, 2020.
- 766 Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog,
767 M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang,
768 Omar Fawzi, et al. Mathematical discoveries from program search with large language models.
769 *Nature*, 625(7995):468–475, 2024.
- 770
- 771 Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro,
772 Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can
773 teach themselves to use tools. In *Thirty-seventh Conference on Neural Information Processing*
774 *Systems*, 2023. URL <https://openreview.net/forum?id=Yacmpz84TH>.
- 775
- 776 Jurgen Schmidhuber. Evolutionary principles in self-referential learning. on learning now to learn:
777 The meta-meta-meta...-hook. Diploma thesis, Technische Universitat Munchen, Germany, 14
778 May 1987. URL <http://www.idsia.ch/~juergen/diploma.html>.
- 779
- 780 Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si,
781 Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, et al. The prompt report: A system-
782 atic survey of prompting techniques. *arXiv preprint arXiv:2406.06608*, 2024.
- 783
- 784 Xuan Shen, Yaohua Wang, Ming Lin, Yilun Huang, Hao Tang, Xiuyu Sun, and Yanzhi Wang. Deep-
785 mad: Mathematical architecture design for deep convolutional neural network. In *Proceedings of*
786 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6163–6173, 2023.
- 787
- 788 Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi,
789 Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. Lan-
790 guage models are multilingual chain-of-thought reasoners. In *The Eleventh International Confer-*
791 *ence on Learning Representations*, 2023.
- 792
- 793 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion:
794 Language agents with verbal reinforcement learning. *Advances in Neural Information Processing*
795 *Systems*, 36, 2023.
- 796
- 797 Kenneth O Stanley and Joel Lehman. *Why greatness cannot be planned: The myth of the objective*.
798 Springer, 2015.
- 799
- 800 Kenneth O Stanley, Jeff Clune, Joel Lehman, and Risto Miikkulainen. Designing neural networks
801 through neuroevolution. *Nature Machine Intelligence*, 1(1):24–35, 2019.
- 802
- 803 Richard S. Sutton. The bitter lesson, 2019. URL [http://www.incompleteideas.net/](http://www.incompleteideas.net/IncIdeas/BitterLesson.html)
804 `IncIdeas/BitterLesson.html`.
- 805
- 806 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 807
- 808 Sai Vemprala, Rogerio Bonatti, Arthur Buckler, and Ashish Kapoor. Chatgpt for robotics: De-
809 sign principles and model abilities. Technical Report MSR-TR-2023-8, Microsoft, Febru-
810 ary 2023. URL [https://www.microsoft.com/en-us/research/publication/](https://www.microsoft.com/en-us/research/publication/chatgpt-for-robotics-design-principles-and-model-abilities/)
811 `chatgpt-for-robotics-design-principles-and-model-abilities/`.
- 812
- 813 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan,
814 and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models.
815 *arXiv preprint arXiv: Arxiv-2305.16291*, 2023a.

- 810 Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos,
811 Charles Blundell, Dhharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn.
812 *arXiv preprint arXiv:1611.05763*, 2016.
- 813
814 Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai
815 Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents.
816 *Frontiers of Computer Science*, 18(6):186345, 2024.
- 817 Rui Wang, Joel Lehman, Jeff Clune, and Kenneth O. Stanley. Poet: open-ended coevolution of
818 environments and their optimized solutions. In *Proceedings of the Genetic and Evolutionary
819 Computation Conference, GECCO '19*, pp. 142–151, New York, NY, USA, 2019. Association for
820 Computing Machinery. ISBN 9781450361118. doi: 10.1145/3321707.3321799.
- 821 Rui Wang, Joel Lehman, Aditya Rawal, Jiale Zhi, Yulun Li, Jeffrey Clune, and Kenneth Stanley.
822 Enhanced poet: Open-ended reinforcement learning through unbounded invention of learning
823 challenges and their solutions. In *International conference on machine learning*, pp. 9940–9951.
824 PMLR, 2020.
- 825
826 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha
827 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
828 models. In *The Eleventh International Conference on Learning Representations*, 2023b.
- 829 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
830 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
831 neural information processing systems*, 35:24824–24837, 2022.
- 832
833 Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li,
834 Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen llm applications via multi-
835 agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023.
- 836 Benfeng Xu, An Yang, Junyang Lin, Quan Wang, Chang Zhou, Yongdong Zhang, and Zhendong
837 Mao. Expertprompting: Instructing large language models to be distinguished experts. *arXiv
838 preprint arXiv:2305.14688*, 2023.
- 839 Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun
840 Chen. Large language models as optimizers. In *The Twelfth International Conference on Learning
841 Representations*, 2024. URL <https://openreview.net/forum?id=Bb4VGOWELI>.
- 842
843 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan
844 Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International
845 Conference on Learning Representations*, 2023. URL [https://openreview.net/forum?
846 id=WE_vluYUL-X](https://openreview.net/forum?id=WE_vluYUL-X).
- 847 Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montserrat Gonzalez
848 Arenas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplik, et al. Language
849 to rewards for robotic skill synthesis. In *Conference on Robot Learning*, pp. 374–404. PMLR,
850 2023.
- 851 Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan, Dongsheng Li, and Deqing Yang. Evoa-
852 gent: Towards automatic multi-agent generation via evolutionary algorithms. *arXiv preprint
853 arXiv:2406.14228*, 2024.
- 854
855 Eliezer Yudkowsky et al. Artificial Intelligence as a positive and negative factor in global risk.
856 *Global catastrophic risks*, 1(303):184, 2008.
- 857
858 Matei Zaharia, Omar Khattab, Lingjiao Chen, Jared Quincy Davis, Heather Miller, Chris Potts,
859 James Zou, Michael Carbin, Jonathan Frankle, Naveen Rao, and Ali Ghodsi. The shift from
860 models to compound ai systems. [https://bair.berkeley.edu/blog/2024/02/18/
861 compound-ai-systems/](https://bair.berkeley.edu/blog/2024/02/18/compound-ai-systems/), 2024.
- 862 Jenny Zhang, Joel Lehman, Kenneth Stanley, and Jeff Clune. OMNI: Open-endedness via mod-
863 els of human notions of interestingness. In *The Twelfth International Conference on Learning
Representations*, 2024a. URL <https://openreview.net/forum?id=AgM3MzT99c>.

864 Shaokun Zhang, Jieyu Zhang, Jiale Liu, Linxin Song, Chi Wang, Ranjay Krishna, and Qingyun
865 Wu. Offline training of language model agents with functions as learnable weights. In *Forty-first*
866 *International Conference on Machine Learning*, 2024b.

867 Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong,
868 and Ji-Rong Wen. A survey on the memory mechanism of large language model based agents.
869 *arXiv preprint arXiv:2404.13501*, 2024c.

870
871 Huaixiu Steven Zheng, Swaroop Mishra, Xinyun Chen, Heng-Tze Cheng, Ed H Chi, Quoc V Le,
872 and Denny Zhou. Take a step back: Evoking reasoning via abstraction in large language models.
873 *arXiv preprint arXiv:2310.06117*, 2023.

874 Pei Zhou, Jay Pujara, Xiang Ren, Xinyun Chen, Heng-Tze Cheng, Quoc V Le, Ed H Chi, Denny
875 Zhou, Swaroop Mishra, and Huaixiu Steven Zheng. Self-discover: Large language models self-
876 compose reasoning structures. *arXiv preprint arXiv:2402.03620*, 2024a.

877
878 Wangchunshu Zhou, Yixin Ou, Shengwei Ding, Long Li, Jialong Wu, Tiannan Wang, Jiamin Chen,
879 Shuai Wang, Xiaohua Xu, Ningyu Zhang, et al. Symbolic learning enables self-evolving agents.
880 *arXiv preprint arXiv:2406.18532*, 2024b.

881 Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen
882 Schmidhuber. Gptswarm: Language agents as optimizable graphs. In *Forty-first International*
883 *Conference on Machine Learning*, 2024.

884
885 Luisa Zintgraf, Sebastian Schulze, Cong Lu, Leo Feng, Maximilian Igl, Kyriacos Shiarlis, Yarin
886 Gal, Katja Hofmann, and Shimon Whiteson. Varibad: Variational bayes-adaptive deep rl via
887 meta-learning. *Journal of Machine Learning Research*, 22(289):1–39, 2021a.

888
889 Luisa M Zintgraf, Leo Feng, Cong Lu, Maximilian Igl, Kristian Hartikainen, Katja Hofmann, and
890 Shimon Whiteson. Exploration in approximate hyper-state space for meta reinforcement learning.
891 In *International Conference on Machine Learning*, pp. 12991–13001. PMLR, 2021b.

892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

SUPPLEMENTARY MATERIAL

TABLE OF CONTENTS

A More Related Work and Future Work	19
A.1 More Related Work	19
A.2 More Future Work	19
B Generalization and Transferability	20
C Prompts	21
D Framework Code	23
E Experiment Details for ARC Challenge	26
F Experiment Details for Reasoning and Problem-Solving Domains	29
G Baselines	30
H Example Agents	31
I Pseudocode of the Meta Agent Search	33
J Impact of Initialization	33
K Cost of Experiments	34

972 A MORE RELATED WORK AND FUTURE WORK

973 A.1 MORE RELATED WORK

974 **AI-Generating Algorithms and AutoML.** Research in AI-Generating Algorithms (AI-
 975 GAs, Clune (2019)) and AutoML (Hutter et al., 2019) aims to replace handcrafted components in AI
 976 systems by learning them. This field has three key pillars: (1) meta-learning architectures, (2) meta-
 977 learning learning algorithms, and (3) generating learning environments and training data (Clune,
 978 2019). Neural Architecture Search (Elsken et al., 2019; Lu et al., 2019; Hu et al., 2021) exem-
 979 plifies the first pillar by automating neural network design, while works like MAML (Finn et al.,
 980 2017) and Meta-RL (Wang et al., 2016; Duan et al., 2017; Norman & Clune, 2023; Zintgraf et al.,
 981 2021a;b) exemplify the second pillar, focusing on “learning to learn” for improved sample efficiency
 982 and generalizability. The third pillar includes works like POET (Wang et al., 2019; Dharna et al.,
 983 2022; Wang et al., 2020) and OMNI-EPIC (Faldor et al., 2024), which generate learning environ-
 984 ments in an open-ended manner. We position Automated Design of Agentic Systems in both the
 985 first and second pillars: meta-learning agentic architectures and leveraging in-context learning to
 986 “learn to learn,” as shown in the ARC challenge (Section 4.1). Furthermore, recent AI-GA and
 987 AutoML advances have also integrated Foundation Models (FMs) to write code, as seen in Fun-
 988 Search (Romera-Paredes et al., 2024) and EoH (Liu et al., 2024), where FMs discover optimization
 989 algorithms. In DiscoPOP (Lu et al., 2024a), FMs program loss functions for preference learning, and
 990 Eureka (Ma et al., 2023) and language-to-reward (Yu et al., 2023) enable FMs to write reward func-
 991 tions for reinforcement learning. OMNI-EPIC (Faldor et al., 2024) allows FMs to create robotics
 992 learning environments. Similarly, we enable FMs to program new agents in code.
 993

994 A.2 MORE FUTURE WORK

- 995 • **More complex domains.** Currently, we only evaluate Meta Agent Search on single-step QA
 996 tasks in this paper. It would be interesting to extend the method to more complex domains, such
 997 as real-world applications involving multi-step interaction with complex environments.
- 998 • **Seeding ADAS with more existing building blocks.** Although we can theoretically allow any
 999 components in agentic systems to be programmed from scratch in the code space, it is not efficient
 1000 in practice. Therefore, it would be interesting to explore ADAS by standing on the shoulders of
 1001 existing human efforts, such as search engine tools, RAG (Lewis et al., 2020), or functions from
 1002 existing agent frameworks like LangChain (LangChainAI, 2022). Additionally, it is interesting
 1003 to support multi-modal capabilities (e.g. vision) in FMs or allow different FMs to be available in
 1004 agentic systems. This will enable the meta agent to choose from different FMs flexibly according
 1005 to the difficulty of the instruction and whether data privacy is a priority.
- 1006 • **Novelty search algorithms.** In Meta Agent Search, the design of the search algorithm is rela-
 1007 tively simple, focusing solely on exploring interesting new designs. A more careful design of the
 1008 search algorithm can be a promising future direction. For example, one could incorporate more
 1009 sophisticated ideas from Quality-Diversity (Mouret & Clune, 2015; Cully & Demiris, 2017), AI-
 1010 generating (Clune, 2019), and Open-ended Algorithms (Faldor et al., 2024; Zhang et al., 2024a;
 1011 Stanley & Lehman, 2015; Stanley et al., 2019). One could also include more classic approaches
 1012 to balance exploration and exploitation (Sutton & Barto, 2018; Liu et al., 2024).
- 1013 • **More Intelligent Evaluation Functions.** In this work, we simply evaluate discovered agents
 1014 on the evaluation set and use the numerical performance results. However, this approach is both
 1015 expensive and misses a lot of information. A promising future direction is to enable the meta
 1016 agent to analyze detailed running logs during the evaluation, which contain rich information on
 1017 the failure and success modes for better debugging and improving agentic systems (Zhou et al.,
 1018 2024b). Also, many tasks involve subjective answer evaluations (Chiang et al., 2024; Lu et al.,
 1019 2024b) that do not have ground-truth answers. It is also important to design novel evaluation
 1020 functions in ADAS to address these tasks. Finally, in this work, we targeted only one domain
 1021 during the search. It would be interesting to explore whether ADAS algorithms can design even
 1022 better generalist agents when specifically searching for agents capable of performing well across
 1023 multiple domains.
- 1024 • **Understanding the emergence of complexity from human organizations.** Beyond potentially
 1025 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 (Hong et al., 2023) or simulating a human town with agents (Park et al., 2023). 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.

B GENERALIZATION AND TRANSFERABILITY

In this section, we present more details of the experiments in Section 4.3 and the complete results of transferring agents across different domains.

For the results shown in Table 3, we use “gpt-4o-2024-05-13” for GPT-4, “claude-3-haiku-20240307” for Claude-Haiku, and “claude-3-5-sonnet-20240620” for Claude-Sonnet.

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 (Wei et al., 2022)	28.0 ± 3.1	34.9 ± 3.2	15.0 ± 2.5	77.8 ± 2.8	88.9 ± 2.2
COT-SC (Wang et al., 2023b)	28.2 ± 3.1	37.8 ± 3.4	15.5 ± 2.5	78.2 ± 2.8	89.0 ± 2.1
Self-Refine (Madaan et al., 2024)	27.5 ± 3.1	38.9 ± 3.4	15.1 ± 2.4	78.5 ± 2.8	89.2 ± 2.2
LLM Debate (Du et al., 2023)	39.0 ± 3.4	43.6 ± 3.4	17.4 ± 2.6	76.0 ± 3.0	88.9 ± 2.2
Step-back Abstraction (Zheng et al., 2023)	31.1 ± 3.2	31.5 ± 3.3	12.2 ± 2.3	76.1 ± 3.0	87.8 ± 2.3
Quality-Diversity (Lu et al., 2024c)	23.8 ± 3.0	28.0 ± 3.1	14.1 ± 2.4	69.8 ± 3.2	80.1 ± 2.8
Role Assignment (Xu et al., 2023)	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

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	Accuracy (%)	
	Math	Reading Comprehension	Multi-task	Science
Manually Designed Agents				
Chain-of-Thought (Wei et al., 2022)	28.0 ± 3.1	64.2 ± 0.9	65.4 ± 3.3	29.2 ± 3.1
COT-SC (Wang et al., 2023b)	28.2 ± 3.1	64.4 ± 0.8	65.9 ± 3.2	30.5 ± 3.2
Self-Refine (Madaan et al., 2024)	27.5 ± 3.1	59.2 ± 0.9	63.5 ± 3.4	31.6 ± 3.2
LLM Debate (Du et al., 2023)	39.0 ± 3.4	60.6 ± 0.9	65.6 ± 3.3	31.4 ± 3.2
Step-back Abstraction (Zheng et al., 2023)	31.1 ± 3.2	60.4 ± 1.0	65.1 ± 3.3	26.9 ± 3.0
Quality-Diversity (Lu et al., 2024c)	23.8 ± 3.0	61.8 ± 0.9	65.1 ± 3.1	30.2 ± 3.1
Role Assignment (Xu et al., 2023)	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

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 (Cobbe et al., 2021), GSM-Hard (Gao et al., 2023), SVAMP (Patel et al., 2021), and ASDiv (Miao et al., 2020) and (2) three domains beyond math adopted in Section 4.2. 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.

C 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.

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 E and F and the framework code is available in Appendix D. 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. These formatting prompts are inspired by Lu et al. (2024a).

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

1134 the implementation steps. The second key (“name”) corresponds to the name of your next agent
 1135 architecture. Finally, the last key (“code”) corresponds to the exact “forward()” function in Python code
 1136 that you would like to try. You must write COMPLETE CODE in “code”: Your code will be part of the
 1137 entire project, so please implement complete, reliable, reusable code snippets.
 1138

1139 Here is an example of the output format for the next agent:
 1140 {“thought”: “**Insights:** Your insights on what should be the next interesting agent. **Overall
 1141 Idea:** your reasoning and the overall concept behind the agent design. **Implementation:** describe
 1142 the implementation step by step.”,
 1143 “name”: “Name of your proposed agent”,
 1144 “code”: “def forward(self, taskInfo): # Your code here”}

1145 ## WRONG Implementation examples:
 1146 [Examples of potential mistakes the meta agent may make in implementation]

1147
 1148
 1149 After the first response from the meta agent, we perform two rounds of self-reflection to make the
 1150 generated agent novel and error-free (Shinn et al., 2023; Madaan et al., 2024).
 1151

1152 Prompt for self-reflection round 1.

1153 [Generated Agent from Previous Iteration]

1154 Carefully review the proposed new architecture and reflect on the following points:
 1155

- 1156
- 1157 1. **Interestingness**: Assess whether your proposed architecture is interesting or innovative
 1158 compared to existing methods in the archive. If you determine that the proposed architecture is not
 1159 interesting, suggest a new architecture that addresses these shortcomings.
 1160 - Make sure to check the difference between the proposed architecture and previous attempts.
 1161 - Compare the proposal and the architectures in the archive CAREFULLY, including their actual
 1162 differences in the implementation.
 1163 - Decide whether the current architecture is innovative.
 1164 - USE CRITICAL THINKING!
 - 1165 2. **Implementation Mistakes**: Identify any mistakes you may have made in the implementation.
 1166 Review the code carefully, debug any issues you find, and provide a corrected version. REMEMBER
 1167 checking “## WRONG Implementation examples” in the prompt.
 - 1168 3. **Improvement**: Based on the proposed architecture, suggest improvements in the detailed
 1169 implementation that could increase its performance or effectiveness. In this step, focus on refining and
 1170 optimizing the existing implementation without altering the overall design framework, except if you
 1171 want to propose a different architecture if the current is not interesting.
 1172 - Observe carefully about whether the implementation is actually doing what it is supposed to do.
 1173 - Check if there is redundant code or unnecessary steps in the implementation. Replace them with
 1174 effective implementation.
 1175 - Try to avoid the implementation being too similar to the previous agent.

1176 And then, you need to improve or revise the implementation, or implement the new proposed
 1177 architecture based on the reflection.

1178 Your response should be organized as follows:

1179 “reflection”: Provide your thoughts on the interestingness of the architecture, identify any mistakes in
 1180 the implementation, and suggest improvements.
 1181 “thought”: Revise your previous proposal or propose a new architecture if necessary, using the same
 1182 format as the example response.
 1183 “name”: Provide a name for the revised or new architecture. (Don’t put words like “new” or “improved”
 1184 in the name.)
 1185 “code”: Provide the corrected code or an improved implementation. Make sure you actually implement
 1186 your fix and improvement in this code.
 1187

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242

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”.

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

```

```

1242 28 \"""
1243 29 Base class for an FM module.
1244 30
1245 31 Attributes:
1246 32 - output_fields (list): Fields expected in the output.
1247 33 - name (str): Name of the FM module.
1248 34 - role (str): Role description for the FM module.
1249 35 - model (str): Model to be used.
1250 36 - temperature (float): Sampling temperature.
1251 37 - id (str): Unique identifier for the FM module instance.
1252 38 \"""
1253 39
1254 40 def __init__(self, output_fields: list, name: str, role='helpful
1255 41 assistant', model='gpt-3.5-turbo-0125', temperature=0.5) -> None:
1256 42 ...
1257 43
1258 44 def generate_prompt(self, input_infos, instruction) -> str:
1259 45 \"""
1260 46 Generates a prompt for the FM.
1261 47
1262 48 Args:
1263 49 - input_infos (list): List of input information.
1264 50 - instruction (str): Instruction for the task.
1265 51
1266 52 Returns:
1267 53 - tuple: System prompt and user prompt.
1268 54
1269 55 An example of generated prompt:
1270 56 ""
1271 57 You are a helpful assistant.
1272 58
1273 59 # Output Format:
1274 60 Reply EXACTLY with the following JSON format.
1275 61 ...
1276 62
1277 63 # Your Task:
1278 64 You will given some number of paired example inputs and outputs.
1279 65 The outputs ...
1280 66
1281 67 ### thinking #1 by Chain-of-Thought hkFo (yourself):
1282 68 ...
1283 69
1284 70 # Instruction:
1285 71 Please think step by step and then solve the task by writing the
1286 72 code.
1287 73 ""
1288 74 \"""
1289 75 ...
1290 76 return system_prompt, prompt
1291 77
1292 78 def query(self, input_infos: list, instruction, iteration_idx=-1) ->
1293 79 list[Info]:
1294 80 \"""
1295 81 Queries the FM with provided input information and instruction.
1296 82
1297 83 Args:
1298 84 - input_infos (list): List of input information.
1299 85 - instruction (str): Instruction for the task.
1300 86 - iteration_idx (int): Iteration index for the task.
1301 87
1302 88 Returns:
1303 89 - output_infos (list[Info]): Output information.
1304 90 \"""
1305 91 ...
1306 92 return output_infos

```



```

1296
129789 def __repr__(self):
129890     return f"{self.agent_name} {self.id}"
129991
129992
130093 def __call__(self, input_infos: list, instruction, iteration_idx=-1):
130194     return self.query(input_infos, instruction, iteration_idx=
130295         iteration_idx)
130396
130396 class AgentSystem:
130497     def forward(self, taskInfo) -> Union[Info, str]:
130598         \"""
130699         Placeholder method for processing task information.
1307100
1307101         Args:
1308102         - taskInfo (Info): Task information.
1309103
1310104         Returns:
1311105         - Answer (Union[Info, str]): Your FINAL Answer. Return either a
1312106           namedtuple Info or a string for the answer.
1313107         \"""
1314     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

```

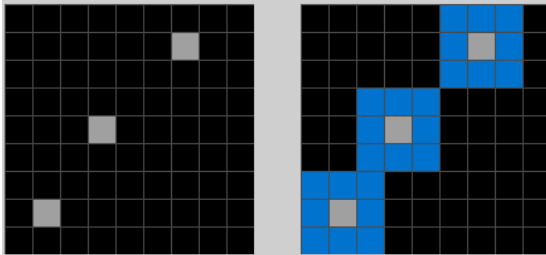
1319 1 def forward(self, taskInfo):
1320 2     # Instruction for initial reasoning
1321 3     cot_initial_instruction = "Please think step by step and then solve
1322 4         the task."
1323 5
1323 5     # Instruction for reflecting on previous attempts and feedback to
1324 6     improve
1325 6     cot_reflect_instruction = "Given previous attempts and feedback,
1326 7         carefully consider where you could go wrong in your latest
1327 8         attempt. Using insights from previous attempts, try to solve the
1328 9         task better."
1329 8     cot_module = FM_Module(['thinking', 'answer'], 'Chain-of-Thought')
1330 9
1330 9     # Instruction for providing feedback and correcting the answer
1331 10    critic_instruction = "Please review the answer above and criticize on
1332 11        where might be wrong. If you are absolutely sure it is correct,
1333 12        output 'True' in 'correct'."
1334 12    critic_module = FM_Module(['feedback', 'correct'], 'Critic')
1335 13
1335 13    N_max = 5 # Maximum number of attempts
1336 14
1336 14    # Initial attempt
1337 15    cot_inputs = [taskInfo]
1338 16    thinking, answer = cot_module(cot_inputs, cot_initial_instruction, 0)
1339 17
1339 18    for i in range(N_max):
1340 19        # Get feedback and correct status from the critic
1341 20        feedback, correct = critic_module([taskInfo, thinking, answer],
1342 21            critic_instruction, i)
1343 22        if correct.content == 'True':
1344 23            break
1345 24
1345 24        # Add feedback to the inputs for the next iteration
1346 25        cot_inputs.extend([thinking, answer, feedback])
1347 26
1347 27        # Reflect on previous attempts and refine the answer
1348 28        thinking, answer = cot_module(cot_inputs, cot_reflect_instruction
1349 29            , i + 1)

```

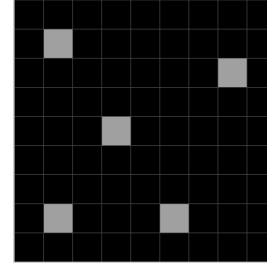
```
1350 | return answer
```

E EXPERIMENT DETAILS FOR ARC CHALLENGE

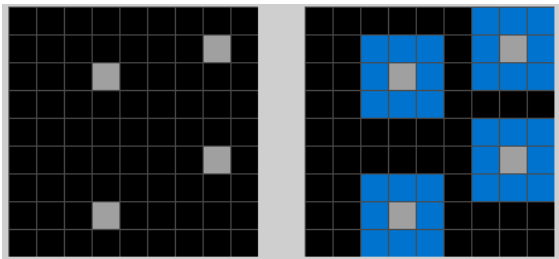
Example Input-output grid #1



Test grid



Example Input-output grid #2



Answer

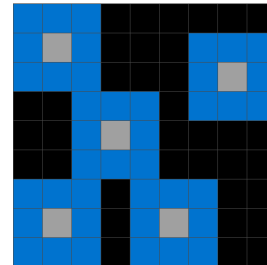


Figure 4: **An example task from the ARC challenge (Chollet, 2019).** Given the input-output grid examples, the AI system is asked to learn the transformation rules and then apply these learned rules to the test grid to predict the final answer.

An example task from the ARC challenge is shown in Figure 4. In the ARC challenge experiments (Section 4.1), 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” (OpenAI, 2024), while discovered agents and baselines are evaluated using “gpt-3.5-turbo-0125” (OpenAI, 2022) 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

1404 outputs, there is also one test input without a known output.
 1405 The inputs and outputs are each “grids”. A grid is a rectangular matrix of integers between 0 and 9
 1406 (inclusive). Each number corresponds to a color. 0 is black.
 1407 Your task is to determine the transformation rule from examples and find out the answer, involving
 1408 determining the size of the output grid for the test and correctly filling each cell of the grid with the
 1409 appropriate color or number.
 1410
 1411 The transformation only needs to be unambiguous and applicable to the example inputs and the test
 1412 input. It doesn’t need to work for all possible inputs. Observe the examples carefully, imagine the grid
 1413 visually, and try to find the pattern.
 1414 **## Examples:**
 1415 **### Example 0:**
 1416 input = [[0,0,0,5,0,0,0,0], [0,0,0,5,0,0,0,0], [0,0,0,4,5,0,0,0], [0,0,0,4,5,4,4,0,0], [0,0,3,3,5,0,0,0,0],
 1417 [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], [0,0,0,5,0,0,0,0], [0,0,0,5,0,0,0,0]]
 1418 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],
 1419 [0,0,0,0]]
 1420 **### Example 1:**
 1421 input = [[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],
 1422 [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]]
 1423 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],
 1424 [0,0,0,0]]
 1425 **### Example 2:**
 1426 input = [[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],
 1427 [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]]
 1428 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],
 1429 [0,0,0,0]]
 1430 **### Test Problem:**
 1431 input = [[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],
 1432 [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]]
 1433 Analyze the transformation rules based on the provided Examples and determine what the output should
 1434 be for the Test Problem.
 1435

1436 Here we present the best agent on ARC discovered by Meta Agent Search.
 1437

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

```

1439 1 # Structured Feedback and Ensemble Agent
1440 2 def forward(self, taskInfo):
1441 3     # Step 1: Generate initial candidate solutions using multiple FM
1442 4     Modules
1443 5     initial_instruction = 'Please think step by step and then solve the
1444 6     task by writing the code.'
1445 7     num_candidates = 5 # Number of initial candidates
1446 8     initial_module = [FM_Module(['thinking', 'code'], 'Initial Solution',
1447 9     temperature=0.8) for _ in range(num_candidates)]
1448 10
1449 11     initial_solutions = []
1450 12     for i in range(num_candidates):
1451 13         thoughts = initial_module[i]([taskInfo], initial_instruction)
1452 14         thinking, code = thoughts[0], thoughts[1]
1453 15         feedback, correct_examples, wrong_examples = self.
1454 16         run_examples_and_get_feedback(code)
1455 17         if len(correct_examples) > 0: # Only consider solutions that
1456 18         passed at least one example
1457 19             initial_solutions.append({'thinking': thinking, 'code': code,
1458 20             'feedback': feedback, 'correct_count': len(
1459 21             correct_examples)})
1460 22
1461 23     # Step 2: Simulate human-like feedback for each candidate solution

```

```

1458 17 human_like_feedback_module = FM_Module(['thinking', 'feedback'], '
1459      Human-like Feedback', temperature=0.5)
1460 18 human_feedback_instruction = 'Please provide human-like feedback for
1461      the code, focusing on common mistakes, heuristic corrections, and
1462      best practices.'
1463 19
1464 20 for sol in initial_solutions:
1465 21     thoughts = human_like_feedback_module([taskInfo, sol['thinking'],
1466      sol['code']], human_feedback_instruction)
1467 22     human_thinking, human_feedback = thoughts[0], thoughts[1]
1468 23     sol['human_feedback'] = human_feedback
1469 24
1470 25 # Step 3: Assign expert advisors to evaluate and provide targeted
1471      feedback
1472 26 expert_roles = ['Efficiency Expert', 'Readability Expert', '
1473      Simplicity Expert']
1474 27 expert_advisors = [FM_Module(['thinking', 'feedback'], role,
1475      temperature=0.6) for role in expert_roles]
1476 28 expert_instruction = 'Please evaluate the given code and provide
1477      targeted feedback for improvement.'
1478 29
1479 30 for sol in initial_solutions:
1480 31     sol_feedback = {}
1481 32     for advisor in expert_advisors:
1482 33         thoughts = advisor([taskInfo, sol['thinking'], sol['code']],
1483         expert_instruction)
1484 34         thinking, feedback = thoughts[0], thoughts[1]
1485 35         sol_feedback[advisor.role] = feedback
1486 36     sol['expert_feedback'] = sol_feedback
1487 37
1488 38 # Step 4: Parse and structure the feedback to avoid redundancy and
1489      refine the solutions iteratively
1490 39 max_refinement_iterations = 3
1491 40 refinement_module = FM_Module(['thinking', 'code'], 'Refinement
1492      Module', temperature=0.5)
1493 41 refined_solutions = []
1494 42
1495 43 for sol in initial_solutions:
1496 44     for i in range(max_refinement_iterations):
1497 45         combined_feedback = sol['feedback'].content + sol['
1498         human_feedback'].content + ''.join([fb.content for fb in
1499         sol['expert_feedback'].values()])
1500 46         structured_feedback = ' '.join(set(combined_feedback.split()
1501         ) # Avoid redundancy
1502         refinement_instruction = 'Using the structured feedback,
1503         refine the solution to improve its performance.'
1504 47         thoughts = refinement_module([taskInfo, sol['thinking'], sol[
1505         'code'], Info('feedback', 'Structured Feedback',
1506         structured_feedback, i)], refinement_instruction, i)
1507 48         refinement_thinking, refined_code = thoughts[0], thoughts[1]
1508 49         feedback, correct_examples, wrong_examples = self.
1509         run_examples_and_get_feedback(refined_code)
1510 50         if len(correct_examples) > 0:
1511 51             sol.update({'thinking': refinement_thinking, 'code':
1512             refined_code, 'feedback': feedback, 'correct_count':
1513             len(correct_examples)})
1514 52             refined_solutions.append(sol)
1515 53
1516 54 # Step 5: Select the best-performing solutions and make a final
1517      decision using an ensemble approach
1518 55 sorted_solutions = sorted(refined_solutions, key=lambda x: x['
1519      correct_count'], reverse=True)
1520 56 top_solutions = sorted_solutions[:3] # Select the top 3 solutions
1521 57
1522 58

```

```

1512 final_decision_instruction = 'Given all the above solutions, reason
1513 over them carefully and provide a final answer by writing the
1514 code.'
1515 final_decision_module = refinement_module(['thinking', 'code'], '
1516 Final Decision Module', temperature=0.1)
1517 final_inputs = [taskInfo] + [item for solution in top_solutions for
1518 item in [solution['thinking'], solution['code'], solution['
1519 feedback']]]
1520 final_thoughts = final_decision_module(final_inputs,
1521 final_decision_instruction)
1522 final_thinking, final_code = final_thoughts[0], final_thoughts[1]
1523 answer = self.get_test_output_from_code(final_code)
1524 return answer

```

F EXPERIMENT DETAILS FOR REASONING AND PROBLEM-SOLVING DOMAINS

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 (OpenAI, 2023). The meta agent uses “gpt-4o-2024-05-13” (OpenAI, 2024), while discovered agents and baselines are evaluated using “gpt-3.5-turbo-0125” (OpenAI, 2022) 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?

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

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 $\implies \Delta E = 3.3 \cdot 10^{-7}$ eV. And $\Delta t = 10^{-11}$ s gives $\Delta E = 3.3 \cdot 10^{-8}$ 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

G BASELINES

In this paper, we implement five state-of-the-art hand-designed agent baselines for experiments on ARC (Section 4.1): (1) Chain-of-Thought (COT) (Wei et al., 2022), (2) Self-Consistency with Chain-of-Thought (COT-SC) (Wang et al., 2023b), (3) Self-Refine (Madaan et al., 2024; Shinn et al., 2023), (4) LLM-Debate (Du et al., 2023), and (5) Quality-Diversity, a simplified version of Intelligent Go-Explore (Lu et al., 2024c).

In addition to these baselines, we implement two more for experiments on Reasoning and Problem-Solving domains (Section 4.2): (6) Step-back Abstraction (Zheng et al., 2023) and (7) Role Assignment (Xu et al., 2023). An example implementation of Self-Refine with our simple framework is shown in Appendix D.

1620 In COT, we prompt the FM to think step by step before answering the question. In COT-SC, we
 1621 sample $N = 5$ answers and then perform an ensemble using either majority voting or an FM query.
 1622 In Self-Refine, we allow up to five refinement iterations, with an early stop if the critic deems the
 1623 answer correct. In LLM-Debate, each debate module is assigned a unique role, such as Physics
 1624 Expert or Chemistry Expert, and the debate lasts for two rounds. In Quality-Diversity, we conduct
 1625 three iterations to collect diverse answers based on previously proposed ones. In Role Assignment,
 1626 we use an FM query to first choose a role from a predefined set, and then use another FM query to
 1627 answer the question by acting within the chosen role.

1628

1629 H EXAMPLE AGENTS

1630

1631 In this section, we present the detailed implementation of three example discovered agents by Meta
 1632 Agent Search shown in Figure 1. The “Multi-Step Peer Review Agent” and “Divide and Conquer
 1633 Agent” were discovered during the search in the Reading Comprehension domain (GPQA) (Rein
 1634 et al., 2023), while the “Verified Multimodal Agent” was discovered during the search in the Math
 1635 domain (MGSM) (Shi et al., 2023).

1636

1637 Code 4: Example discovered agent: Multi-Step Peer Review Agent

```

1638 1 def forward(self, taskInfo):
1639 2     initial_instruction = "Please think step by step and then solve the
1640 3     task."
1641 4     critique_instruction = "Please review the answer above and provide
1642 5     feedback on where it might be wrong. If you are absolutely sure
1643 6     it is correct, output 'True' in 'correct'."
1644 7     refine_instruction = "Given previous attempts and feedback, carefully
1645 8     consider where you could go wrong in your latest attempt. Using
1646 9     insights from previous attempts, try to solve the task better."
1647 10    final_decision_instruction = "Given all the above thinking and
1648 11    answers, reason over them carefully and provide a final answer."
1649 12
1650 13    FM_modules = [FM_module(['thinking', 'answer'], 'FM Module', role=
1651 14    role) for role in ['Physics Expert', 'Chemistry Expert', 'Biology
1652 15    Expert', 'Science Generalist']]
1653 16    critic_modules = [FM_module(['feedback', 'correct'], 'Critic', role=
1654 17    role) for role in ['Physics Critic', 'Chemistry Critic', 'Biology
1655 18    Critic', 'General Critic']]
1656 19    final_decision_module = FM_module(['thinking', 'answer'], 'Final
1657 20    Decision', temperature=0.1)
1658 21
1659 22    all_thinking = [[] for _ in range(len(FM_modules))]
1660 23    all_answer = [[] for _ in range(len(FM_modules))]
1661 24    all_feedback = [[] for _ in range(len(FM_modules))]
1662 25
1663 26    for i in range(len(FM_modules)):
1664 27        thinking, answer = FM_modules[i](taskInfo, initial_instruction)
1665 28        all_thinking[i].append(thinking)
1666 29        all_answer[i].append(answer)
1667 30
1668 31    for i in range(len(FM_modules)):
1669 32        for j in range(len(FM_modules)):
1670 33            if i != j:
1671 34                feedback, correct = critic_modules[j](taskInfo,
1672 35                all_thinking[i][0], all_answer[i][0]),
1673 36                critique_instruction)
1674 37                all_feedback[i].append(feedback)
1675 38
1676 39    for i in range(len(FM_modules)):
1677 40        refine_inputs = [taskInfo, all_thinking[i][0], all_answer[i][0]]
1678 41        + all_feedback[i]
1679 42        thinking, answer = FM_modules[i](refine_inputs,
1680 43        refine_instruction)
1681 44        all_thinking[i].append(thinking)
1682 45
1683 46

```

```

1674     all_answer[i].append(answer)
1675
1676     final_inputs = [taskInfo] + [all_thinking[i][1] for i in range(len(
1677         FM_modules))] + [all_answer[i][1] for i in range(len(FM_modules))
1678         ]
1679     thinking, answer = final_decision_module(final_inputs,
1680         final_decision_instruction)
1681
1682     return answer

```

Code 5: Example discovered agent: Divide and Conquer Agent

```

1684 1 def forward(self, taskInfo):
1685 2     # Step 1: Decompose the problem into sub-problems
1686 3     decomposition_instruction = "Please decompose the problem into
1687 4     smaller, manageable sub-problems. List each sub-problem clearly."
1688 5     decomposition_module = FM_Module(['thinking', 'sub_problems'], '
1689 6     Decomposition Module')
1690 7
1691 8     # Step 2: Assign each sub-problem to a specialized expert
1692 9     sub_problem_instruction = "Please think step by step and then solve
1693 10    the sub-problem."
1694 11    specialized_experts = [FM_Module(['thinking', 'sub_solution'], '
1695 12    Specialized Expert', role=role) for role in ['Physics Expert', '
1696 13    Chemistry Expert', 'Biology Expert', 'General Expert']]
1697 14
1698 15    # Step 3: Integrate the sub-problem solutions into the final answer
1699 16    integration_instruction = "Given the solutions to the sub-problems,
1700 17    integrate them to provide a final answer to the original problem.
1701 18    "
1702 19    integration_module = FM_Module(['thinking', 'answer'], 'Integration
1703 20    Module', temperature=0.1)
1704 21
1705 22    # Decompose the problem
1706 23    thinking, sub_problems = decomposition_module([taskInfo],
1707 24    decomposition_instruction)
1708 25
1709 26    # Ensure sub_problems is a string and split into individual sub-
1710 27    problems
1711 28    sub_problems_list = sub_problems.content.split('\n') if isinstance(
1712 29    sub_problems.content, str) else []
1713 30
1714 31    # Solve each sub-problem
1715 32    sub_solutions = []
1716 33    for i, sub_problem in enumerate(sub_problems_list):
1717 34        sub_problem_info = Info('sub_problem', decomposition_module.
1718 35        __repr__(), sub_problem, i)
1719 36        sub_thinking, sub_solution = specialized_experts[i % len(
1720 37        specialized_experts)]([sub_problem_info],
1721 38        sub_problem_instruction)
1722 39        sub_solutions.append(sub_solution)
1723 40
1724 41    # Integrate the sub-problem solutions
1725 42    integration_inputs = [taskInfo] + sub_solutions
1726 43    thinking, answer = integration_module(integration_inputs,
1727 44    integration_instruction)
1728
1729    return answer

```

Code 6: Example discovered agent: Verified Multimodal Agent

```

1725 1 def forward(self, taskInfo):
1726 2     # Instruction for generating visual representation of the problem
1727 3     visual_instruction = "Please create a visual representation (e.g.,
    diagram, graph) of the given problem."

```



```

1728 4
1729 5 # Instruction for verifying the visual representation
1730 6 verification_instruction = "Please verify the accuracy and relevance
1731 7 of the visual representation. Provide feedback and suggestions
1732 8 for improvement if necessary."
1733 9
1734 10 # Instruction for solving the problem using the verified visual aid
1735 11 cot_instruction = "Using the provided visual representation, think
1736 12 step by step and solve the problem."
1737 13
1738 14 # Instantiate the visual representation module, verification module,
1739 15 and Chain-of-Thought module
1740 16 visual_module = FM_Module(['visual'], 'Visual Representation Module')
1741 17 verification_module = FM_Module(['feedback', 'verified_visual'], '
1742 18 Verification Module')
1743 19 cot_module = FM_Module(['thinking', 'answer'], 'Chain-of-Thought
1744 20 Module')
1745 21
1746 22 # Generate the visual representation of the problem
1747 23 visual_output = visual_module([taskInfo], visual_instruction)
1748 24 visual_representation = visual_output[0] # Using Info object
1749 25 directly
1750
1751 # Verify the visual representation
1752 feedback, verified_visual = verification_module([taskInfo,
1753 visual_representation], verification_instruction)
1754
1755 # Use the verified visual representation to solve the problem
1756 thinking, answer = cot_module([taskInfo, verified_visual],
1757 cot_instruction)
1758 return answer

```

I PSEUDOCODE OF THE META AGENT SEARCH

In this section, we provide the pseudocode for the Meta Agent Search algorithm to clarify its implementation and workflow. The pseudocode outlines the iterative process of designing, evaluating, and refining agents using a meta agent, as described in the main text.

Algorithm 1 Meta Agent Search Algorithm

```

1764 1: Input: Target domain validation data, maximum iterations  $N$ 
1765 2: Output: Archive of discovered agents
1766 3: Initialize archive  $\mathcal{A}$  with baseline agents (e.g., Chain-of-Thought, Self-Refine)
1767 4: for  $i = 1$  to  $N$  do
1768 5:   Design Step: Meta agent generates a new agent:
1769 6:     (a) Outputs design reasoning
1770 7:     (b) Implements the design in code
1771 8:     (c) Performs two self-reflection steps to ensure novelty and correctness
1772 9:   Evaluation Step: Evaluate the new agent on target domain validation data:
1773 10:    (a) If the agent produces errors during evaluation, refine the design up to 5 iterations
1774 11:    (b) Re-run the evaluation after each refinement
1775 12:   Update Step: Add the refined agent and its evaluation metrics to the archive  $\mathcal{A}$ 
1776 13: end for
1777 14: Return: Final archive  $\mathcal{A}$ 

```

J IMPACT OF INITIALIZATION

One of the key claims of our work is that the *code space* representation allows for better utilization of existing human efforts (Section 2), enabling a more efficient search process than starting entirely

from scratch. To further investigate the effects of initialization, we conducted experiments where the Meta Agent Search algorithm was run without any initial agent designs, contrasting with our standard approach that incorporates human-designed solutions into the search process.

The results, presented in Table 6, demonstrate that even without initial agent designs, Meta Agent Search discovers agents that outperform all hand-crafted baselines across all evaluated domains. This finding underscores the robustness of our method, as it effectively leverages the inherent structure of the code space to explore and optimize agent designs.

Interestingly, while the inclusion of good initial solutions generally leads to improved performance, the math domain exhibited a unique outcome: starting from scratch resulted in superior performance. We hypothesize that the absence of predefined design patterns in this case encouraged a broader and more diverse exploration of reasoning strategies within the limited number of iterations. Such diversity appears particularly beneficial for math tasks, which demand flexible and varied approaches to reasoning.

This observation opens up an intriguing avenue for future research: exploring how the choice and quality of initialization impact search effectiveness across different domains. For instance, it would be valuable to identify conditions under which starting without initial solutions may yield performance gains, or to design strategies that combine the advantages of both initialization and broad exploration.

Table 6: **Performance comparison of Meta Agent Search with and without initial agent designs across multiple domains.** The results show that even without initialization, Meta Agent Search outperforms hand-designed baselines in all domains. However, incorporating initial solutions generally leads to better performance, except in the math domain, where starting without initialization yields superior results.

Agent Name	F1 Score		Accuracy (%)	
	Reading Comprehension	Math	Multi-task	Science
State-of-the-art Hand-designed Agents				
Chain-of-Thought (Wei et al., 2022)	64.2 ± 0.9	28.0 ± 3.1	65.4 ± 3.3	29.2 ± 3.1
COT-SC (Wang et al., 2023b)	64.4 ± 0.8	28.2 ± 3.1	65.9 ± 3.2	30.5 ± 3.2
Self-Refine (Madaan et al., 2024)	59.2 ± 0.9	27.5 ± 3.1	63.5 ± 3.4	31.6 ± 3.2
LLM Debate (Du et al., 2023)	60.6 ± 0.9	39.0 ± 3.4	65.6 ± 3.3	31.4 ± 3.2
Step-back Abstraction (Zheng et al., 2023)	60.4 ± 1.0	31.1 ± 3.2	65.1 ± 3.3	26.9 ± 3.0
Quality-Diversity (Lu et al., 2024c)	61.8 ± 0.9	23.8 ± 3.0	65.1 ± 3.3	30.2 ± 3.1
Role Assignment (Xu et al., 2023)	65.8 ± 0.9	30.1 ± 3.2	64.5 ± 3.3	31.1 ± 3.1
Automated Design of Agentic Systems on Different Domains				
Meta Agent Search (Empty Initialization)	73.9 ± 0.9	67.5 ± 3.3	68.5 ± 3.3	32.7 ± 3.2
Meta Agent Search	79.4 ± 0.8	53.4 ± 3.5	69.6 ± 3.2	34.6 ± 3.2

K COST OF EXPERIMENTS

A single run of search and evaluation on ARC (Section 4.1) costs approximately \$500 USD in OpenAI API costs, while a run within the reasoning and problem-solving domains (Section 4.2) 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 6, 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.