# SelfAI: Building a Self-Training AI System with LLM Agents

Xiao Wu[1,2], Ting-Zhu Huang[1*], Liang-Jian Deng[1], Xiaobing Yu[5],
Yu Zhong[1], Shangqi Deng[3], Ufaq Khan[2], Jianghao Wu[6],
Xiaofeng Liu[4], Imran Razzak[2], Xiaojun Chang[2], Yutong Xie[2*]

[1]University of Electronic Science and Technology of China, [2]Mohamed
bin Zayed University of Artificial Intelligence, [3]Xian Jiaotong University,
[4]Yale University, [5]Washington University in St. Louis, [6]Monash
University .

*Corresponding author(s). E-mail(s): tingzhuhuang@126.com;
yutong.xie@mbzuai.ac.ae;

## Abstract

Recent work on autonomous scientific discovery has leveraged LLM-based agents to integrate problem specification, experiment planning, and execution into end-to-end systems. However, these frameworks are often confined to narrow application domains, offer limited real-time interaction with researchers, and lack principled mechanisms for determining when to halt exploration, resulting in inefficiencies, reproducibility challenges, and under-utilized human expertise. To address these gaps, we propose *SelfAI*, a general multi-agent platform that combines a User Agent for translating high-level research objectives into standardized experimental configurations, a Cognitive Agent powered by LLMs with optimal stopping criteria to iteratively refine hyperparameter searches, and an Experiment Manager responsible for orchestrating parallel, fault-tolerant training workflows across heterogeneous hardware while maintaining a structured knowledge base for continuous feedback. We further introduce two novel evaluation metrics, Score and $AUP_D$, to quantify discovery efficiency and search diversity. Across regression, NLP, computer vision, scientific computing, medical imaging, and drug discovery benchmarks, SelfAI consistently achieves strong performance and reduces redundant trials compared to classical Bayesian optimization and LLM-based baselines, while enabling seamless interaction with human researchers.

# 1 Main

In recent years, large language models (LLMs) [1–3] have fundamentally reshaped the landscape of AI-driven research. Advances in reasoning [4, 5], multimodal understanding [6, 7], and autonomous tool use [8–13] have positioned LLMs as central components for accelerating scientific workflows.

Early LLM-based scientific research demonstrated that LLMs can extract actionable scientific knowledge to guide experiments [14–17] or answer professional questions [18–20]. As reasoning and tool-interaction capabilities matured, research expanded to cross-stage planning [21–23], hypothesis and idea generation [24–27], multimodal knowledge integration [28], and conducting experiments [17, 29]. Meanwhile, LLMs have driven the development of tasks such as reaction prediction [30–32], material discovery [18, 33–35], chemical synthesis design and molecular property optimization [29, 30, 36, 37], and biomedicine research [17, 38–41]. Building on these advances, a new generation of scientific discovery systems has emerged, spanning AI scientific assistants across the full spectrum of scientific workflows [16, 18, 42, 43], automated scientific discovery systems (ASDSs) specialized for code generation and experiment design [29, 36, 44], and increasingly capable research agents that incorporate automated debugging, iterative idea refinement, and scientific writing [26, 45, 46]. Along this evolutionary trajectory, LLM-centered ASDS have begun to integrate inference, execution, and feedback mechanisms into unified closed-loop pipelines, achieving varying degrees of autonomy across scientific domains.

Despite rapid progress, existing ASDS frameworks still exhibit fundamental limitations. Many systems [24, 44, 45, 47] (such as AIRA and Scientist-V2) excel at translating research intent into executable code or experimental procedures, thereby validating the "executability" of ASDS and primarily relying on final success rate (medal rate) or best performance as metrics. More autonomous platforms like MLGym [48] and MLAgentBench [49] enable system-level control, allowing agents to execute end-to-end benchmarks, collect results, and conduct standardized evaluations. Moreover, while recent LLM-driven optimization methods [50–53] such as LLM4EO showcase that LLMs can infer evolutionary tendencies and synthesize more effective operators, they operate at the level of local modification rules rather than scientific reasoning in trajectories. Overall, these approaches do not explicitly enhance scientific reasoning during experimentation, such as identifying optimal stopping points, discovering efficient trial trajectories, or assessing the structural quality of reasoning.

Building reasoning-centric ASDS therefore requires a shift from execution-oriented autonomy toward cognitive autonomy: systems must not only carry out experiments but also analyze, reflect, and revise their reasoning throughout the experimental process. However, current ASDS systems provide little support for these cognitive-layer behaviors. The absence of mechanisms for trajectory-level reasoning assessment and LLM-driven decision-making limits both efficiency and accuracy, especially in settings with large search spaces, stochastic outcomes, or complex interdependencies among hypotheses, resulting in inefficient exploration, elevated computational costs, and limited real-time adaptability. These limitations highlight a fundamental need for ASDS architectures that incorporate reasoning-driven search, adaptive learning mechanisms, and structural evaluation of scientific trajectories.
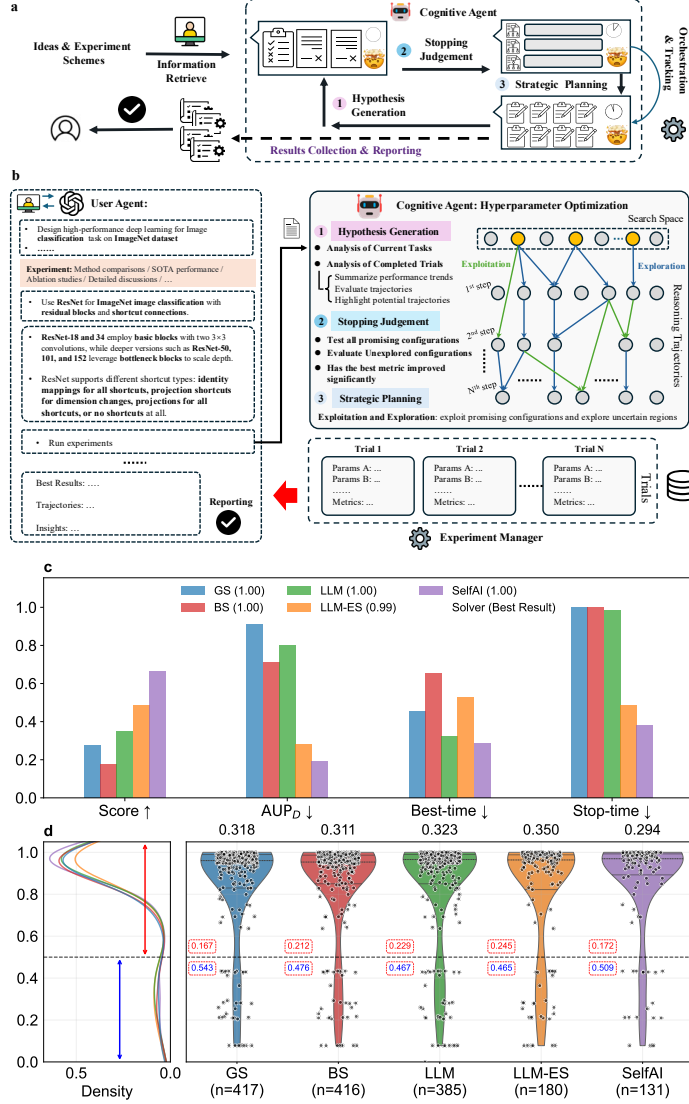
**Fig. 1**: **SelfAI Framework for Automated Scientific Experimentation. a**, Holistic architecture of the multi-agent system, which transforms various experiments in the research process into a structured workflow. **b**, User intentions, comprising ideas and experiment schemes, are transformed into structured configurations via a predefined prompt. These inputs are processed through successive stages: hypothesis generation, strategic planning, trial execution, and result collection. **c**, Performance distribution across 11 tasks demonstrates the framework's ability, when powered by GPT-4o-mini, to prioritize high-performance regions without sacrificing exploration. **d**, Trial counts for each solver shown in **c**, accompanied by quantile lines, density distributions, and performance variability across the global and two evaluation regions. Higher values in low-performance regions promote rapid escape, while lower values in high-performance regions enable localized refinement.

To address these issues, we propose **SelfAI**, a unified multi-agent self-training pipeline for continuous adaptability and transparent collaboration. As shown in Fig. 1, SelfAI forms a closed loop that integrates user intent, cognitive reasoning, and experimental orchestration. Upstream, a **User Agent** converts high-level objectives and exploratory questions into standardized experiment configurations. These configurations enable the **Cognitive Agent** to analyze performance metrics, conduct reasoning over historical outcomes, and adjust the search trajectory, driving gradual improvement [54, 55]. A resource-aware **Experiment Manager** ensures efficient execution by handling resource scheduling, environment provisioning, adaptive checkpointing, and comprehensive experiment logging. During high-level parameter optimization, the Cognitive Agent and Experiment Manager jointly estimate the potential benefit of continued training and apply an optimal-stopping criterion to terminate unpromising trials. To evaluate scientific reasoning quality, we introduce two complementary evaluation metrics, Score (Discovery Efficiency) and $\text{AUP}_D$ (Area Under the Performance Diversity). Score aggregates, across tasks, the normalized improvement over the search space together with penalties for discovering good configurations late and for stopping far from the best-found point. $\text{AUP}_D$ explicitly encodes how broadly a solver explores the search space by summarizing the performance-diversity tradeoffs across the entire trajectory, enabling detailed analysis of exploration behavior and stopping decisions in long-term searches. Together, these metrics enable quantitative assessment of exploration behavior, reasoning structure, and stopping decisions in long-term autonomous experiments.

In summary, SelfAI provides a continuously self-training system that autonomously iterates through cycles of intent-centric understanding, trajectory-aware reasoning, and adaptive strategic planning, enabling sustained performance gains with minimal human intervention. Through checkpoint tolerance, zero-code parallelization, and dynamic resource management, SelfAI substantially reduces redundant computation, improves resource efficiency, and accelerates search convergence in long-horizon scientific workflows. In addition, SelfAI integrates seamlessly with modern MLOps pipelines [56–58], facilitating collaborative experiment coordination, end-to-end traceability of model evolution, and robust privacy protection for sensitive configurations and user data. Collectively, these capabilities make SelfAI a practical, scalable, and operationally mature platform for rapid scientific discovery. The key contributions of this paper are given as follows:

1. **Autonomous Scientific Discovery System:** We design a cohesive framework of specialized LLM agents, including the *User Agent*, *Cognitive Agent Framework*, and *Experiment Manager*, which collaboratively translate high-level research objectives into iteratively refined experimental workflows, including trajectory analysis, hypothesis generation, and best stopping decision-making.
2. **Novel Evaluation Metrics and Optimal Stopping:** We propose two new metrics, *Score* (discovery efficiency) and $\text{AUP}_D$ (area under the performance–diversity curve), to jointly quantify the efficiency and diversity of hyperparameter exploration. An embedded optimal stopping criterion automatically terminates
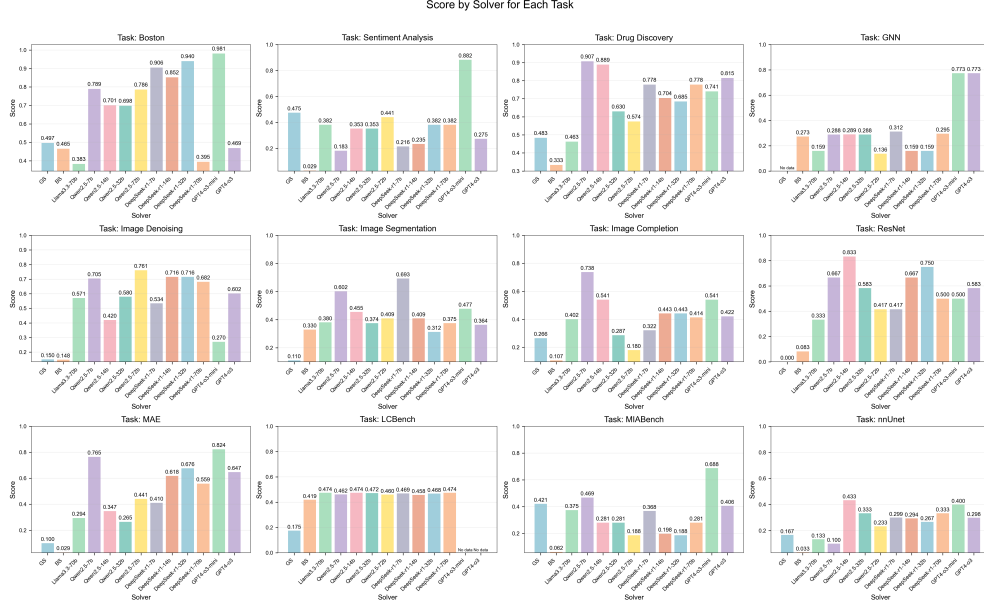
**Fig. 2**: Scores among all solvers across different tasks to measure the best stopping criterion.

    unpromising trials, effectively balancing exploration and exploitation while minimizing computational overhead.

3. **Comprehensive Validation of LLM-based Reasoning in Scientific Discovery:** We evaluate SelfAI across diverse tasks, including regression, sentiment analysis, computer vision, medical imaging, and drug discovery. The results show that SelfAI effectively improves upon excellent baseline models in both results exploration and performance.

## 2 Performance of SelfAI on Benchmark Data

Scientific discovery is driven by diverse physical environments and user intents. To effectively reason across diverse scientific discovery tasks, our approach integrates multiple LLM agents and a suite of tools: User Agent, Cognitive Framework, and Experiment Management, addressing the complex reasoning and planning problem of scientific discovery. In our experiments, we evaluate the SelfAI across 6 domains and 12 tasks. The compared methods contain different suites of LLMs, such as OpenAI-o3 [1], Llama3.3 [59], Qwen2.5 [2], and DeepSeek-R1 [3]. We used grid search and the Tree Parzen Estimator (TPE) optimizer [60, 61] (referred to as BS), which applies a modified Bayesian inference method using a spiral search. All agents are tested on four Nvidia A100 GPUs. A public implementation of SelfAI will be available on the project website. We also provide reasoning processes for reproducibility.
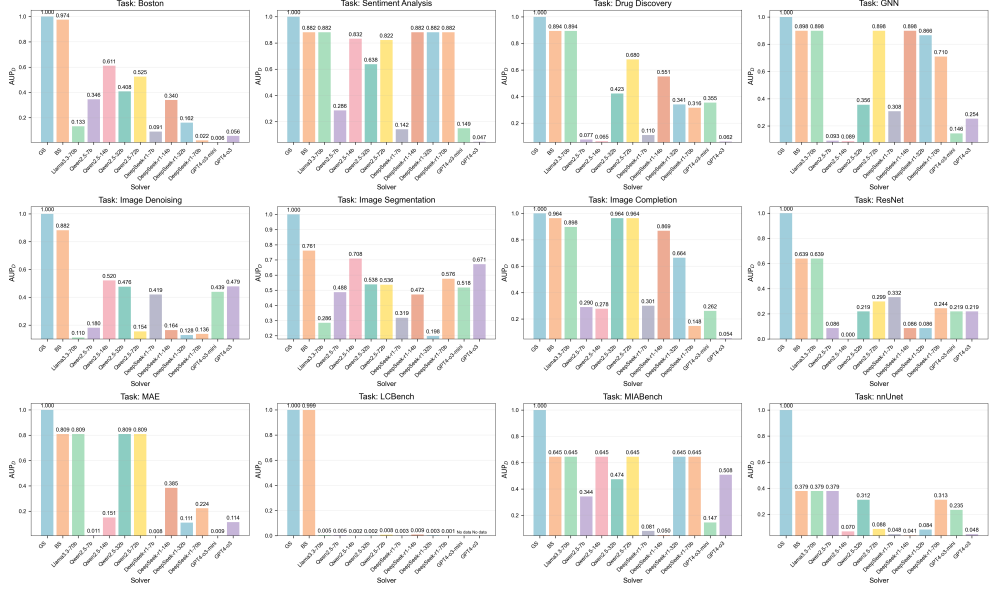
**Fig. 3**: Diverse Metrics ($\text{AUP}_D$) among all solvers across different tasks to evaluate trajectory diversity.

**Metrics.** We introduce four metrics: Score, $\text{AUP}_D$, $t_{\text{best}}$, $t_{\text{stop}}$. Score measures a comprehensive search performance that combines optimality and stopping efficiency. $\text{AUP}_D$ quantifies the diversity of explored high-quality solutions. $t_{\text{best}}$ represents the normalized time to first find the best result, while $t_{\text{stop}}$ indicates when the search terminates. Specifically, the reasoning process that finds the best result but delays stopping exhibits low $t_{\text{best}}$ and high $t_{\text{stop}}$. Conversely, late discovery with immediate stopping yields high values for both. An ideal solver should minimize both times. In addition, a high $\text{AUP}_D$ reflects broad exploration, whereas a low value indicates focused search, which may also result from rapid convergence near the starting point. Thus, the most valuable optimization strategy achieves a high Score with a low $\text{AUP}_D$, enabling fast discovery of optimal configurations with minimal wasted exploration. More details are illustrated in the Supplementary Section A.4.

**Benchmarks.** We collect benchmarks composed of multiple tasks of 6 primary categories, designed to simulate challenging scenarios in scientific discovery. The benchmark requires the solver to engage in scientific reasoning by analyzing experiments, interpreting observations, and adapting strategies, uncovering novel patterns, functional relationships, or optimal experimental configurations. Compared with hyperparameter search, these processes reflect the inductive reasoning essential to real-world scientific inquiry. Therefore, this benchmark assesses the model's reasoning ability, not just its numerical optimization capabilities. All benchmark data is collected through actual experimental runs, and some results are reported in published literature, spanning both discrete and continuous hyperparameter search spaces, from low-dimensional to high-dimensional settings (see Appendix B for more details).

Existing methods [24, 45, 49] explore final performance or medal rates using reward functions and prompts specifically designed for innovative idea generation, code generation, and pass@K strategies. In the line of LLM-driven optimization, LLM4EO [50] uses LLMs exclusively to synthesize and evolve evolutionary operators that determine which parts of a candidate solution should be modified during the search. The LLM is used to generate and update operators to shape the local search behavior of the evolutionary algorithm. These implementations do not focus on scientific reasoning across trials during the exploration process, making direct comparisons with our benchmark infeasible. Consequently, we used the readily accessible GPT4-o3-mini model and compared SelfAI with LLM and LLM-ES methods on 11 small/medium-sized tasks across different domains in Fig. 1c. The LLM and LLM-ES methods are fine-tuned by using the structured prompts in Supplementary Section D. The LLM solver uses a search prompt that recommends new trials based on prior results, while LLM-ES augments this with an early-stopping module that decides whether to terminate exploration. Fig. 1c presents that our SelfAI outperforms LLM and LLM-ES in all four metrics, validating SelfAI's effectiveness in integrating domain knowledge to accelerate scientific discovery. Furthermore, we focus on the reasoning capability of the original LLM model because it is difficult to be fair when comparing different methods. Fig. 1d shows that SelfAI promptly focuses on exploiting higher-performance regions while maintaining a balanced distribution for effective exploration. Figs. 2 and 3 summarize the benchmark performance based on these two metrics. The benchmarks show that SelfAI consistently achieves excellent results across various domains and tasks for different LLM models. Detailed comparisons covering Score, $AUP_D$, Best-Time, Stop-Time, and Hit Rate are provided in Supplementary Table B2.

**Failure cases.** Several key failure modes of LLM-based solvers are shown in Supplementary Section B.1. First, the hit rate for identifying the best result varies substantially across models, with some solvers (particularly DeepSeek-R1 series) failing to reach the global optimum due to premature stopping. Second, limited context windows cause incomplete reasoning over long trajectories, preventing solvers from fully incorporating earlier experimental signals. Third, computational fragility leads to non-monotonic performance under seemingly minor perturbations in reasoning strategy. These findings indicate that while SelfAI improves overall efficiency, reasoning-driven exploration remains sensitive to solver stability and stopping behavior.

## 2.1 Machine Learning

**Boston house pricing prediction.** We evaluate SelfAI on the Boston housing price prediction task [62] using a random forest regression model. A total of 162 trials were conducted across a search space defined by five hyperparameters: i.e., n-estimators $= [100, 200, 300]$, max-depth $= [None, 10, 20]$, min-samples-split $= [2, 5, 10]$, min-samples-leaf$=[1,2,5]$ and max-features $= [\text{"sqrt"},\text{"log2"}])$. As shown in Figs. 2 and 3 (and Supplementary Table C3), GPT-4o-mini achieves by far the highest Score (0.9811) and the 1st ranking, indicating outstanding optimization efficiency. Within the DeepSeek-R1 family, the 32B and 7B variants secure 2nd and 3rd places respectively, demonstrating the best balance between rapid convergence (low Best-Time)

and reasonable search diversity among open-source models. In contrast, DeepSeek-R1-70B and Llama3.3-70B exhibited consistent exploration but lacked effective stopping criteria, leading to longer convergence times. Notably, nearly all models achieved an identical and high Best Result ( 0.841), indicating that while most can find a correct solution, they differ substantially in optimization efficiency and reliability.

**Sentiment analysis.** We perform experiments for sentiment analysis that focus on identifying opinions, emotions, and attitudes expressed in text. All experiments are found in [63] with standard experimental settings, where the pre-trained Word2Vec embeddings [64] are used as input features. An LSTM network is then employed to model sentence sequences, converting them into a multi-class classification problem. This setup serves to evaluate the agent's reasoning and optimization capabilities in textual domains. As shown in Supplementary Table C4, GPT-4o-mini outperforms all other LLMs, achieving the highest Score (0.8824) by discovering the optimal configuration extremely early while maintaining reasonable search diversity. DeepSeek-R1 series and Qwen2.5 series exhibit significantly lower Scores, with many failing to converge quickly or occasionally missing the global optimum. The largest frontier model GPT-4o also underperforms markedly (Score 0.2745, rank 11), suggesting that scale alone is insufficient for effective iterative hyperparameter reasoning in recurrent neural architectures.

## 2.2 Scientific Computing for Image Completion

For scientific computing fields, we selected the tensor decomposition method [65]. The method is an important tool for high-dimensional data analysis and is crucial in applications such as data compression [66], computational acceleration [67], and multi-modal data fusion [28]. All LLMs are provided with identical mathematical knowledge about TW decomposition generated by GPT-4o. As shown in Supplementary Table C5, Qwen2.5 (7b and 14b) and GPT4-o3-mini perform well and rank highly. Qwen2.5-72b, on the other hand, may overly rely on general reasoning capabilities and fail to balance mathematical knowledge and reasoning capabilities. In the DeepSeek series, although the DeepSeek-R1-7b model can't find the optimal solution, its search strategy was acceptable with a moderate Score, demonstrating its ability to perform mathematical reasoning. DeepSeek-R1-70b and GPT-4o, despite excellent search diversity and early stopping, received a mediocre score, reflecting that their exploration strategy was not well aligned with the optimization goal of tensor decomposition, possibly failing to find the optimal performance between data compression and computational efficiency.

## 2.3 Computer Vision

**A. SIREN** We employed SIREN (Sinusoidal Representation Networks) [68] to evaluate our framework on image segmentation and denoising tasks. Leveraging sine activations and coordinate-based inputs, SIREN excels at representing high-frequency signals through continuous implicit representations, making it widely used in scientific computing and physics-based problems. However, its performance is highly sensitive to hyperparameters (e.g., learning rate and regularization strength, etc.), requiring careful tuning per dataset. The unsupervised nature of SIREN further increases the

risk of training instability or divergence with improper settings, making it a strong test case for SelfAI's hyperparameter optimization capability.

Fig. 4 illustrates hyperparameter search trajectories for image segmentation using SIREN (see also Supplementary Figs. B4, B5 and B6), where the surface, smoothed from the original steep, multi-peak data, reveals distinct search behaviors. Given the same three initial points, the tree-structured Bayesian optimizer (BS) follows a spiral trajectory that underutilizes promising starting regions and fails to explore broadly, often converging to local minima. In contrast, LLM-based optimizers infer trends from evaluated points, incorporate causal reasoning, and explore more broadly to efficiently locate the global optimum. They also monitor progress and halt when improvements plateau, a capability absent in traditional methods.

Quantitative results (Supplementary Tables C6–C7) reveal distinct task-dependent performance profiles. In segmentation, DeepSeek-R1-7B achieves the highest Score (0.693) and ranks first, followed by Qwen2.5-7B and GPT-4-o3-mini. Larger models such as Qwen2.5-72B and Llama-3.3-70B rank only mid-range, indicating that scale alone does not ensure effective optimization. In denoising, the landscape shifts: Qwen2.5-72B delivers the best performance (Score 0.761), while DeepSeek-R1-14B and DeepSeek-R1-32B jointly occupy second place. Notably, models that perform well in segmentation tasks do not necessarily perform well in denoising tasks, suggesting that the effectiveness of the solver largely depends on the degree of matching between its inference strategy and the task.

**B. Image Classification** We also benchmark two commonly used supervised learning methods in computer vision: Masked Autoencoder (MAE) [69] and ResNet [70]. These represent fundamentally different learning paradigms where hyperparameter optimization plays a crucial role in achieving the latest performance.

**Mask Autoencoder (MAE)** MAE introduces two key hyperparameters: training strategies (linear detection vs. fine-tuning) and masking rate, which are explored over a range: [0.10, 0.90] with an interval of 0.1. The masking rate directly affects the difficulty of the reconstruction task and the quality of the learning representation, while the training strategy determines how to adapt to the downstream task. In our experiments (see Supplementary Table C8 for more details), GPT-4o-mini achieved explicit performance with efficiency, demonstrating the extraordinary ability to identify the best mask configuration and training strategies. Qwen2.5-7B achieves early convergence and disciplined stopping. The DeepSeek-R1 series shows consistent scaling behavior: performance improves from 7B to 14B to 32B, but slightly drops at 70B. In traditional solvers, GS and BS lag significantly behind LLM-based solvers, highlighting the advantages of intelligent hyperparameter optimization methods. Apart from Qwen2.5-14b and DeepSeek-R1-7b, most solvers can find the optimal solution, but significant differences exist among the various methods in terms of search efficiency and convergence speed.

**ResNet Family** On the ImageNet ResNet hyperparameter benchmark, LLMs in the 7B–32B range demonstrate superior optimization efficacy. The top-performing solvers, such as Qwen2.5-14B (1st), DeepSeek-R1-32B (2nd), and Qwen2.5-7B with DeepSeek-R1-14B (tied 3rd), consistently identify the optimal ResNet configuration within one or two trials, halting immediately with near-perfect sample efficiency. In contrast, larger variants from the same model families drop to middle or lower rankings,

consuming over half the search budget with notably higher exploration overhead. While GPT-series models, as large-scale counterparts, still achieve competitive results, the 7B–32B class exhibits a clear advantage in balancing accuracy and efficiency. For the architectural search on ImageNet, we explored key dimensions such as depth and bottleneck design. Qwen2.5-14B attains theoretical peak performance, with its 7B and 32B versions completing further exploration efficiently. These results reflect an ability to identify ideal network configurations with minimal wasted exploration.

**LCBench** We evaluate SelfAI on the standard AutoML benchmark [71] using Bayesian optimization over 2000 rounds of hyperparameter search, covering critical parameters including learning rate, batch size, network depth, and dropout rate. As shown in Supplementary Table C10, DeepSeek-R1-70B, Qwen2.5-14B, and Llama3.3-70B secure the top three positions in the ranking, with DeepSeek-R1-70B attaining the highest overall score and the lowest search diversity. Notably, leading models in the 7B–32B parameter range efficiently identify near-optimal hyperparameter configurations by evaluating only a minimal set of candidate solutions, achieving an order-of-magnitude improvement in stop-time compared to classical methods. These findings demonstrate diminishing returns with scaling, as larger models within the same series (e.g., Qwen2.5-72B) typically underperform comparable models of medium size, while the 7B model can achieve robust performance comparable to models several times its size.

## 2.4 Medical Image Analysis

In medical image analysis, nnU-Net [72] is a landmark framework known for its strong generalization and automated design. Despite its adaptive network structure, preprocessing, and training strategies for various segmentation tasks, the exploration of novel architectures persists. To address this, nnU-Net-Revisited [73] establishes a comprehensive benchmark including 19 mainstream models (CNN-based [74–76], Transformer-based [77, 78], and Mamba-based [79, 80]), offering a solid basis for fair evaluation. Results are shown in Figs. 2-3 and Supplementary Tables C11- C12. On the BTCV dataset [81], GPT4-o3-mini achieves the highest Score (0.6875) and Rank 1, demonstrating strong optimization efficacy, while Qwen2.5-7B secures second place, outperforming all larger variants in its family. In nnU-Net hyperparameter tuning for the BraTS dataset [82], Qwen2.5-14B attains the top ranking (Score 0.4333), followed closely by GPT4-o3-mini, with Qwen2.5-32B and DeepSeek-R1-70B tying for third. Notably, the 7B and 72B models outperformed similar models of medium size (14B–32B), indicating that larger model size does not guarantee better optimization.

## 2.5 Imbalanced Node Classification

Graph neural networks (GNNs), as powerful tools for processing relational data, are playing an increasingly critical role in numerous cutting-edge scientific fields. In biomedicine, GNNs have become a core component of protein structure prediction models such as AlphaFold3 [39], enabling precise modeling of spatial interactions between amino acid residues and ushering in a new era in structural biology. In drug discovery, they are widely used for molecular property prediction [39, 83, 84], compound-target interaction identification [85, 86], and novel drug generation [87]. Similarly,

**Fig. 4**: Illustration of the optimized trajectory for the SIREN method for image segmentation. Green points are suggested points before reaching the optimal points. Red points are redundant suggestions when reaching out to the optimal points and failing to stop trials. The ⋆ is the optimal point. We show the serialization recommendations provided by LLM through the labeled numbers.

GNNs demonstrate exceptional performance in engineering domains, including traffic prediction [88, 89], cybersecurity, and chip design [90, 91]. Given this broad utility and the sensitivity of GNN performance to hyperparameter settings, we evaluate GraphSAGE on the imbalanced Cora benchmark [92] under the setup of [93].

11

Figs. 2 and 3 represent that the GPT4-o3 series delivered near-optimal results, achieving the fastest discovery of optimal solutions and efficient stopping decision while maintaining high search diversity, underscoring its robust optimization capability in scientific computing scenarios. In contrast, Qwen2.5 models, (except 72b), easily delivered consistent but middling results across scales, while larger models generally exhibit late convergence and higher unnecessary exploration, reflecting potential inherent limitations in understanding the structural properties and optimization requirements of specialized domains such as protein interaction networks and molecular graph structures (see Supplementary Table C13 for more results). These results indicate that SelfAI can effectively address challenges in GNN hyperparameter optimization. The framework's strong performance on graph learning tasks suggests its potential as a general-purpose optimizer for scientific domains, including protein engineering and drug development.

## 2.6 Drug Discovery

In drug discovery, accurate bioactivity prediction is vital for early-stage virtual screening and compound prioritization, enabling shorter development cycles and lower preclinical costs. Following the evaluation practices and modeling standards summarized by Korotcov et al. [94], this study employs the Chagas EP20 dataset [95], which measures the activity of compounds against Trypanosoma cruzi, a neglected but medically significant target in tropical disease research. The core challenge lies in learning an effective mapping from molecular representations, such as SMILES sequences [96] or graph-based encodings [97] to experimentally measured bioactivity [98]. We adopt SelfAI to adaptively refine performance, assess search trajectories, and determine principled stopping points for this task.

Experimental results (see Supplementary Table C14) reveal distinct performance patterns among solvers in hyperparameter optimization. Qwen2.5-7B and 14B secure first and second place, respectively, both achieving high Scores with minimal search diversity. GPT-4o ranks third, while DeepSeek-R1-7B and 70B tie for fourth. A key finding is the clear performance drop at larger scales; Qwen2.5-32B, Qwen2.5-72B, and DeepSeek-R1-32B rank in the bottom half, demonstrating that increased parameters do not guarantee better optimization. In contrast, traditional optimization methods (such as GS and BS) ranked poorly in this task, underscoring the clear advantage of large language model-based intelligent optimizers in bioactivity prediction. These results collectively demonstrate that in high-dimensional, noisy biochemical scenarios, language model-driven optimizers offer efficient convergence and a superior balance between exploration and exploitation.

Finally, the ranking heatmap (Fig. B3) and average performance analysis (Supplementary Table B2) reveal several consistent trends. GPT-4o-mini emerges as the most consistently effective solver, strong inductive reasoning, rapid adaptation to emerging evidence, and stable progression across diverse scientific domains. Mid-sized models such as Qwen2.5-7B and Qwen2.5-14B display competitive and robust performance relative to their scale, often exhibiting disciplined evidence-based refinement during exploration. In contrast, larger models (Qwen2.5-72B and DeepSeek-R1-70B) exhibit

higher variance: they perform well on specific tasks but frequently show delayed break-throughs, unstable trajectories, and inconsistent cumulative progress. These patterns suggest that larger models tend to spend more time exploring alternative possibilities rather than leveraging early evidence, resulting in delayed commitment and reduced adaptability during the search process.

# 3 Discussion

The emergence of the autonomous scientific discovery system (ASDS) represents a paradigm shift in scientific research, with the potential to accelerate discovery by delegating structured experimental processes to artificial intelligence. Yet most existing agentic ML frameworks and benchmark-driven ASDS, such as MLE-bench [99], emphasize an agent's ability to generate and execute candidate solutions, offering limited support for adaptive trajectory management, such as reasoning about when exploration should continue, how search strategies should evolve, or why certain regions of the space warrant further investigation. SelfAI addresses this gap by integrating structured intent interpretation, multi-agent reasoning, and principled stopping criteria into a general reasoning framework. Across diverse domains, this design enables the system to allocate computational effort where it is most informative, rather than exhaustively enumerating the search space.

Under this goal, SelfAI is meticulously designed as a general-purpose platform capable of performing high-level reasoning across diverse tasks, thereby accelerating the experimental search process. Evaluations across a comprehensive benchmark (12 tasks across 6 different domains) uncover that SelfAI strategically allocates limited computational resources to the most promising performance spaces or research directions and terminates unproductive trajectories, exhibiting human-like cognitive flexibility in dynamically adjusting the search trajectory. This optimal stopping behavior effectively alleviates the most persistent inefficiencies in current ASDS.

The newly proposed Score and $\text{AUP}_D$ metrics are crucial for quantifying this advantage, revealing SelfAI's remarkable balance between rapidly finding discovery efficiency (Score) and maintaining healthy exploration diversity ($\text{AUP}_D$). These metrics further illuminate the limitations of traditional methods like grid search and Bayesian optimization, which often lack the adaptive reasoning required to escape local optima or terminate unfruitful trials. Our findings also suggest that in scientific domains, effective scientific discovery demands specialized reasoning capabilities beyond mere model scale. Although large-scale models like GPT-4o and DeepSeek-R1-70B occasionally excel in broad exploration, their performance is frequently surpassed by more compact, purpose-driven models such as Qwen2.5-7B and DeepSeek-R1-7B, which demonstrate superior stopping efficiency and search consistency. This indicates that for navigating complex non-convex optimization environments, the structured reasoning and hypothesis generation implemented in SelfAI's cognitive agent may play a more decisive role than model scale alone. The framework demonstrates that purpose-driven architectural design can effectively harness the capabilities of moderately-sized models to deliver robust scientific discovery performance.

The implications of this framework extend far beyond accelerated experimental search and improved benchmark performance. First, the architecture of SelfAI embodies a practical blueprint for human-AI symbiosis in science. The User Agent acts as a seamless translator, allowing human researchers to interact with the system at the level of scientific intent and high-level questions, rather than through low-level code or configuration files. This design formally integrates human expertise as a guiding input, creating a collaborative loop where the AI manages the scale and complexity of the search while the human provides strategic direction and domain insight. This addresses a critical shortcoming of "black-box" automation by fostering reproducibility and trust, as every action of the AI is traceable to a human-defined objective. Furthermore, this closed-loop process allows researchers to shift from using AI for scientific research to exploring how to conduct scientific research optimally and to investigate efficient search strategies and problem-related endpoints under controlled conditions.

SelfAI has the potential to evolve from a tool into a collaborative research partner, enhancing the cognitive and experimental capabilities of human scientists. This evolution would be driven by several key advancements. For instance, the User Agent could be enhanced to incorporate autonomous code generation, thereby broadening the range of experimental configurations without requiring explicit user pre-definition. Furthermore, while the current iteration of SelfAI relies on the static knowledge embedded within its underlying LLM (Large Language Model), future versions could leverage advanced memory management and RAG approaches to alleviate context length limitations and dynamically integrate a richer, more up-to-date body of domain-specific knowledge. Collectively, these envisioned capabilities underscore SelfAI's foundational role in shaping a more adaptive, collaborative, and intelligent paradigm for autonomous scientific discovery.

# 4 The Design of SelfAI for Scientific Discovery

SelfAI is a general scientific discovery system designed to automate and accelerate the end-to-end scientific discovery process. Its core function is to transform users' high-level research intentions (hyperparameter optimization, novel algorithm design, model architecture validation, or ablation studies) into structured, executable workflows. The system operates in an integrated loop, comprising three core agents: a user agent, a cognitive agent, and an experiment manager, which together constitute the overall experimental workflow in Fig. 1.

## 4.1 User Agent

The User Agent, as the user interface of SelfAI, is designed to be interactive, user-controllable, and adaptable to evolving research needs. It interprets high-level scientific intent and translates natural-language objectives into structured, machine-readable experimental configurations. Through iterative clarification, it helps researchers specify objectives, constraints, and search spaces, thereby establishing the experimental environment in which the Cognitive Agent operates. For example, requests such as "achieve state-of-the-art accuracy on CIFAR-10 using a CNN" or "identify the most influential

method for image classification" are reformulated into precise experimental specifications (see Supplementary Material A.1). Unlike static intent-generation prompts, the User Agent supports continuous intervention: researchers may pause execution, inspect reasoning traces, adjust constraints, or redesign experimental protocols at any stage without restarting the workflow.

## 4.2 Cognitive Agent

SelfAI's cognitive agent is central to accelerating scientific discovery. After receiving structured experimental configurations from the user agent and accumulated experimental history from the experiment manager, the agent synthesizes all available information, analyzes experimental trajectories, evaluates observed performance trends, and identifies promising regions in the search space. This achieves a delicate balance between developing these high-potential regions and exploring areas of greater uncertainty, ultimately deriving the next optimal search strategy. A significant feature of our cognitive agent is its introduction of principled stopping criteria, addressing a key limitation of traditional AI systems, which often continue exploring under conditions of diminishing returns. By conducting evidence-based evaluations of research progress, the agent can determine when further experiments on a particular hypothesis are unlikely to yield significant new insights, thereby reallocating resources to more promising research directions (see Supplementary Material A.2 for more details). This capability improves resource utilization efficiency in the automated discovery process.

## 4.3 Experiment Manager

The Experiment Manager executes the experimental plans proposed by the Cognitive Agent, each of which has already been structured as executable training jobs, specifying model architectures, datasets, and hyperparameters. For every job, the manager initializes these specified configurations, launches the training and evaluation pipeline, and records all associated outputs. It tracks runtime status, handles execution failures, and captures performance metrics as well as relevant logs and metadata. These records constitute a continuously updated structured knowledge base, providing a reliable foundation for the cognitive agent to analyze search progress, identify promising regions of the search space, and decide when to stop exploration. Comprehensive execution procedures and implementation details are provided in Supplementary Material A.3.

Beyond executing trials, the Experiment Manager also manages the computational environment in which scientific experiments are carried out. Through the Experiment Manager, resources are allocated reliably, execution workloads are scheduled efficiently, and checkpointing and recovery mechanisms are coordinated seamlessly. Anomalies or failed runs are surfaced to the Cognitive Agent as feedback for complex adaptive reasoning. As a result, the overall research process is conducted within an efficient, reliable, and reasoning-driven environment.

**Table 1**: Comparison of SelfAI with related AI research frameworks and benchmarks across system-level, agent-specific, and task-specific capabilities.

| Functions | Ours | Code LLaMA [53] | MLGym [48] | AI Co-Scientist [45] | AIRA [24] | MLAgentBench [49] | Optuna [61] |
|---|---|---|---|---|---|---|---|
| Interactive Research | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | ✓ |
| Flexible Artifacts | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ | ✗ |
| Privacy and Security | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| Trajectory Analysis | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Hypothesis Generation | ✓ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ |
| Strategic Planning | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| Causal Inference | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| Adaptive Learning | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| Job Scheduling | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| Checkpoint Management | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Experiment Tracking | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓(*) |
| Zero-Code Parallelization | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Hypothesis Optimization | ✓ | ✓(*) | ✗ | ✗ | ✗ | ✗ | ✓ |
| Self-Evaluation | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Benchmark Suite | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ |

Note: The comparison is structured in four blocks: (1) System-level functions, (2) Cognitive functions primarily handled by the Cognitive Agent, (3) Execution functions managed by the Experiment Manager, and (4) Performance in optimization tasks. ✓(*) denotes basic support.

## 4.4 Comparison of AI research systems

We contextualize SelfAI's capabilities within the current landscape of AI-assisted research tools. Table 1 presents a systematic comparison against representative frameworks across multiple dimensions of scientific functionality. Existing ASDSs often specialize in code- or idea-level hypothesis generation under fully autonomous pipelines that primarily optimize medal rates within 24 hours. In contrast, SelfAI integrates intent interpretation, iterative scientific reasoning, and adaptive trajectory optimization into a unified discovery workflow. These capabilities allow the system to refine strategic planning based on accumulated experimental evidence, rather than executing predetermined and static generalization procedures. Overall, these advantages position SelfAI as a general-purpose framework capable of supporting efficient, reliable, and reasoning-driven scientific workflows that extend the scope of existing agent or hyperparameter optimization benchmarks.

# 5 Data availability

All datasets used in this study are available for download at https://github.com/XiaoXiao-Woo/SelfAI.

# 6 Code availability

The SelfAI repository is available as Supplementary Software. Updated versions can be found at https://github.com/XiaoXiao-Woo/SelfAI.

# Appendix A    LLM Agents in SelfAI System for Scientific Discovery

The SelfAI framework supports the complete lifecycle of AI systems, covering design, training, evaluation, and deployment. It integrates three core agents: the User Agent, the Cognitive Agent, and the Experiment Manager, and provides an integrated toolkit for autonomous model training, experiment orchestration, and advanced reasoning, enabling a scalable and adaptive workflow.

## A.1    User Agent: Idea Interaction and Experiment Configuration

In autonomous research systems, the User Agent is typically powered by a general-purpose LLM and acts as the primary interface between human researchers and the system. These systems build an "idea-to-experiment" autonomous research process where user access to detailed experimental workflows is limited. Such approaches can streamline research execution but may also restrict researchers' flexibility in iteratively shaping and adapting experiments. By contrast, our User Agent unifies idea generation and experiment configuration into a continuous, interactive process, empowering researchers to refine initial concepts, explore alternative approaches, and adjust experimental parameters at any stage.

During the idea interaction process, the User Agent understands the task background, dissects method details, and defines precise experimental objectives across diverse scientific domains. At this stage, the User Agent can help users pinpoint innovative research directions and generate field-relevant ideas. Once an idea emerges, the configuration interaction phase leverages prompt templates to clarify user objectives and formulate strategic experimental plans, including initial hypotheses and search configurations. This process converts high-level manual queries into standardized and reproducible experimental configurations, as illustrated in the following examples. We insert the specified template into the prompt and trigger the prompt to fill out the specified configuration file.

Consequently, the User Agent can transform ideas into system and user information. The system information provides the Cognitive Agent with essential content, including role-playing instructions and task-specific knowledge (such as task benchmarking, method evaluation, comparative analysis, ablation studies, and experimental designs). This information equips the Cognitive Agent to interpret high-level research objectives accurately and to execute relevant reasoning and planning steps effectively.

In parallel, the user information focuses on the practical aspects of experimental planning by supplying detailed specifications and keywords, such as search space definitions, the number of trials, optimization criteria, and resource constraints, which are critical for configuring and running experiments. Together, these two types of information enable seamless coordination between conceptual ideation and concrete execution, supporting a workflow in which researchers can iteratively refine their experiments and monitor progress in real time.

**User Agent Prompt for Configuration Interaction**

**System:** ...
**User:**
`{{SUMMARIZED_CONTENT}}`
Please fill out the content following the YAML format:

```yaml
# Template for Configuration Interaction
# to convert Idea Interaction into YAML format
- role: system
  content:
    model: $modelName
    description: You are a $role specializing in
                 studying $taskName. Please provide
                 professional and detailed answers.
    task: $taskName
    basic_idea: $basic_idea
    search_space:
      $hyper_name1: []
      $hyper_name2: []
      ...
    link: $paperLink
    instrustion: Complete instrucions under limited
        trials.
- role: user
  content:
    max_trials: $maxTrials
    trials: []
```

## A.2 Cognitive Agent

In our SelfAI framework, the Cognitive Agent continuously ingests and adapts domain knowledge to tackle complex reasoning problems during experiments, much like a human researcher. At each stage, the agent first conducts a comprehensive analysis of the task and completed experiments to clarify the purpose and key information of the task and generates preliminary hypotheses. Then, it evaluates whether these hypotheses are worth trying through stop judgment to inform strategic planning. It then generates specific schemes to probe unexplored regions of the solution space, creating a non-Markovian chain of evolving plans grounded in prior states.

### A.2.1 Reasoning for Scientific Discovery

LLM agents exhibit remarkable reasoning capabilities, making them well-suited for scientific discovery. As illustrated in Fig. A1, we present an advanced reasoning framework for generating optimal stopping solutions, an approach that can be extended to various scientific discovery tasks.
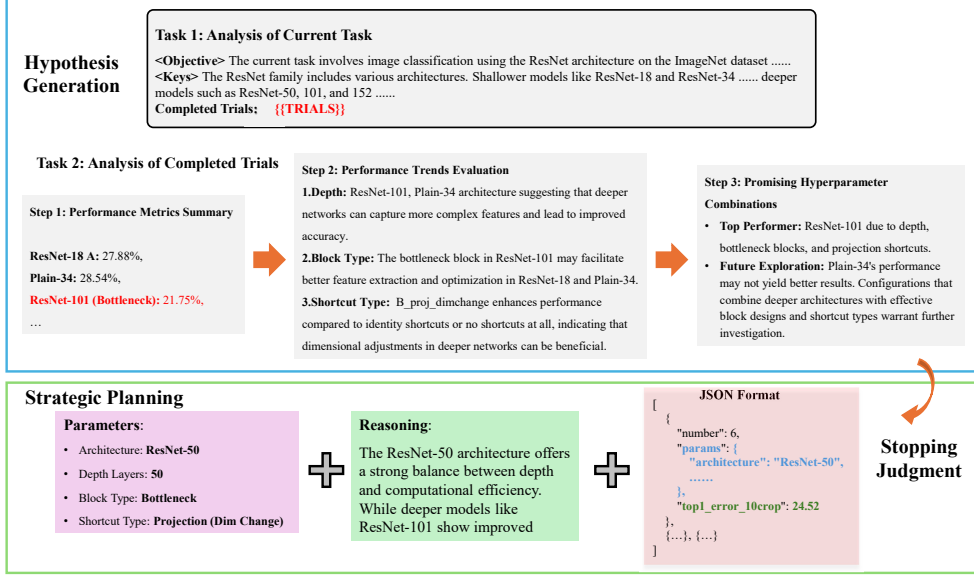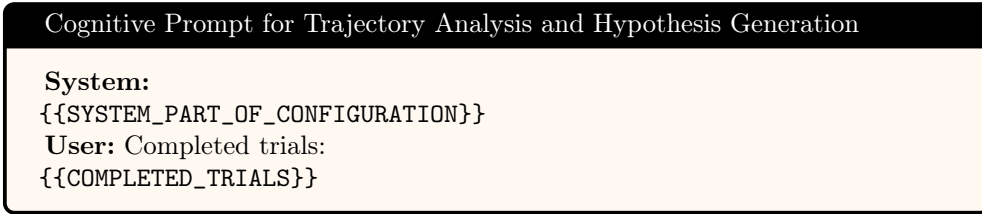
**Fig. A1**: **Illustration of the Cognitive Agent.** The overall reasoning process involves several key steps: Hypothesis Generation (analysis of the current task and completed trials), Stopping Judgment, and Strategic Planning. Strategic Planning develops experimental schemes based on the analyzed hypotheses.

## A. Trajectory Analysis and Hypothesis Generation

Specifically, the agent begins with **Task 1: Analyze the Current Task**, where it meticulously interprets the prompt to establish the task's core objectives, constraints, and hyperparameters. This initial analysis is used to perform a structured trajectory evaluation, creating a coherent chain of thought that guides all subsequent strategy and exploration. This initial analysis is used to perform a structured trajectory evaluation, creating a coherent chain of thought that guides all subsequent reasoning and planning strategies.

Next, in **Task 2: Analysis of Completed Trials**, the agent systematically evaluates previous outcomes to generate new hypotheses. It specifically analyzes the performance trends of individual parameters and their combinations to identify promising directions in the broader hyperparameter space and propose novel configurations. This deep trend analysis allows the agent to anticipate promising regions of the search space and detect valuable reasoning trajectories that might otherwise be overlooked.

---

**Cognitive Prompt for Trajectory Analysis and Hypothesis Generation**

**System:**
`{{SYSTEM_PART_OF_CONFIGURATION}}`
**User:** Completed trials:
`{{COMPLETED_TRIALS}}`

---

> Task 1: Analyze the current task
> Understand current tasks, basic ideas, objectives, and hyperparameters.
> Task 2: Analysis of Completed Trials
> Step 1: Summarize performance metrics for completed trials.
> Step 2: Evaluate performance trends for hyperparameters.
> Step 3: Highlight promising hyperparameter combinations.

**B. Best Stopping Judgement** As described in Sect. **??**, we introduce the optimal stopping criterion to guide the prompt design that balances exploration and exploitation. During Best Stopping Judgement, the cognitive agent first determines whether the current performance metrics surpass those of the initial configuration and conducts an optimal stopping judgment. Building on the prior analysis of completed trials, the agent systematically evaluates stopping conditions to avoid testing all configurations.

By using prompt instructions, the optimal stopping judgment evaluates all completed trials by leveraging insights from trial analysis, observed performance trends, and identified key findings. In practice, we implement this judgment process via a structured prompt template comprising two main tasks and explicit stopping criteria. The language model then analyzes the relationship between completed trials and unexplored areas. The prompt is organized as follows:

> **Cognitive Prompt for Best Stopping Judgement**
>
> **User:** Completed trials:
> `{{COMPLETED_TRIALS}}`
> The following **Search Space** contains **unexplored** trials. `{{TRIALS}}`
> Instructions: Task 1: Review Analysis of Completed Trials (trial analysis, performance trends, highlights, and other insights)
> Task 2: Decide Whether to Stop Optimization
> Based on the above analysis and **Completed Trials**, determine whether the optimization process should be stopped.
> Carefully analyze each of the following stop rules and provide a short (1-2 sentences) justification for whether it is met:
> 1. Have all promising configurations identified based on performance trends been tested?
> 2. Is it unlikely that unexplored configurations will perform better based on the observed trends and the law of diminishing returns?
> 3. Has the best metric improved significantly?
> Step 2: Decide whether **all** conditions are met.
> If **all** criteria in Step 1 are met, Answer: Yes, with confidence score: `{{CONFIDENCE_SOCRE}}`. Otherwise, Answer: No with confidence score: `{{CONFIDENCE_SOCRE}}`.

**C. Strategic Planning** Following the reasoning and stopping judgment phases shown in Fig. A1, the agent reviews all executed trials to synthesize key observations on performance trends, notable configurations, and areas of consistent improvement. Next,

it translates its abstract hypotheses into a concrete experimental plan. This strategic process balances two competing objectives: exploiting promising configurations and exploring uncertain regions. Exploitation focuses on refining or building upon previously successful configurations, whereas exploration targets under-examined areas that show significant potential for performance gains. Each proposed trial is explicitly justified by insights from the preceding reasoning phase, ensuring the search is deliberate and well-founded. Finally, each recommendation is formalized as an intended trial and delegated to the Experiment Manager for execution.

As trials advance and incremental improvements diminish, the Experiment Manager maintains a dynamic record of completed actions and continuously refines the remaining search space. Strategic planning directs the optimization trajectory toward convergence, and the system employs an effective stopping strategy to terminate exploration in a rational and efficient manner.

Notably, the LLM-based agent demonstrates human-like exploratory behavior: it initially focuses on refining hypotheses around high-performing regions and gradually extends attention to areas of greater uncertainty. This pattern reflects a synergy of intelligence-driven reasoning and adaptive planning.

---

**Cognitive Prompt for Strategic Planning**

**User:** Instructions:
Task 1: Review Analysis of All Completed Trials
Completed trials:
`{{COMPLETED_TRIALS}}`
The following **Search Space** contains **unexplored** trials:
 `{{TRIALS}}`
Instructions:
Task 2: Optimization Recommendation
Recommend exactly `{{N_JOBS}}` promising trials from the provided **Search Space** (include both number and params).
**Rules:**
1. "params" MUST include:
 `{{HYPERNAME}}`
2. All selected 'params' must match exactly with the provided **Search Space**. Do NOT leave out any key.
3. Use the analysis in **Task 1** (trial analysis, performance trends, highlights, and other insights) to guide selection.
4. Based on the above analysis, explore under-explored regions only when there is clear evidence of potential performance gain.
5. Do not mix, modify, or create new values.
6. You MUST not output any JSON blocks in this part.
7. You MUST provide reasoning for each recommendation.

## A.3 Experiment Manager

The Experiment Manager is responsible for experiment orchestration and recovery, including resource management, task allocation, and progress tracking. These capabilities enable efficient coordination of multi-instance parallel optimization, maximize resource utilization, and enhance training robustness:

1) **Resource Management**. The Experiment Manager monitors user program resource consumption and dynamically allocates available GPU, TPU, and memory resources. This granular allocation optimizes workload distribution across computing units and ensures stable execution of all trials.

2) **Fault recovery and Checkpoint Reconnection**. In case of system interruptions or suboptimal model performance, the Experiment Manager reports failures to the cognitive agent. The Experiment Manager performs preliminary diagnostics, identifies potential issues, adjusts training parameters, and resumes training from the latest checkpoint.

3) **Multi-Instance Parallel Optimization**. SelfAI instantiates each user program to run across diverse physical environments, independent of the target program framework. The Experiment Manager coordinates multi-instance parallel training, synchronizes execution, and concurrently tests various configurations, thereby shortening overall training time and improving generalization across datasets and model parameters. For each parallel experiment, the Experiment Manager identifies and supplies necessary runtime parameters, ensuring experiments are conducted under the same environment optimized by the cognitive agent.

SelfAI maintains detailed logs of all configurations, training runs, and evaluation metrics, and provides targeted optimization suggestions for each completed experiment. This intelligent decision-making refines experimental design and reduces redundant search in automated scientific discovery. The collaboration between the Experiment Manager and the cognitive agent establishes a self-improving workflow. By conducting advanced reasoning and analyzing trial result patterns through hypothesis-driven exploration, they uncover relationships between settings and model performance, steering the search toward optimal configurations. This ensures that parallel training is not merely a brute-force process, but an intelligent exploration of the training space that adapts based on real-time outcomes.

## A.4 Evaluation for Reasoning Trajectories

While the reasoning process of LLMs in problem-solving often generates a diverse range of discrete insights and multiple potential chains of thought, this diversity, while valuable for exploration and exploitation, can pose a challenge to coherent reasoning evaluation across various experimentation and discovery phases. We propose a systematic evaluation metric that captures both the diversity of reasoning perspectives and the overall coherence of the reasoning process, ensuring a more robust and comprehensive assessment of reasoning capabilities.

**Optimal Stopping Criteria** In this work, we collected the best value point and the stop point from trials. Based on Optimal Stopping Criteria [100, 101], we can define

the following measure formulas,

$$\text{Gain} = \frac{1}{N} \sum_{i=0}^{N-1} \frac{v_i^* - v_{i,\min}}{v_{i,\max} - v_{i,\min}} \tag{A1}$$

$$t_{\text{best}} = \frac{1}{N} \sum_{i=0}^{N-1} t_i^{\text{best}} = \frac{1}{N} \sum_{i=0}^{N-1} \frac{m_i}{M_i} \tag{A2}$$

$$t_{\text{stop}} = \frac{1}{N} \sum_{i=0}^{N-1} t_i^{\text{stop}} = \frac{1}{N} \sum_{i=0}^{N-1} \frac{n_i}{M_i} \tag{A3}$$

and where $N$ is the number of tasks. For the $i$-th task, $M_i$ is the number of completed trials. $m_i$ is the best value point index. $t_i^{\text{best}}$ is the cost of the best value. In addition, we set the stop point index, $n_i$, where $t_i^{\text{stop}}$ is better when $n_i$ is closer to the best value point index. $S_i$ means that the binarized value.

To obtain a comprehensive measure, we combine the last three measures, i.e., Rel, $t_{\text{best}}$, and $t_{\text{stop}}$, where we utilize Rel in underestimated penalty, then $t_{\text{best}}$ and $t_{\text{stop}}$ are the time penalty ($P_{\text{best}}$ and $P_{\text{stop}}$). Thus, the total penalty is

$$P_{\text{total}} = \frac{t_i^{\text{stop}} + t_i^{\text{best}}}{2} \tag{A4}$$

Finally, the score is denoted as

$$\text{Score} = \frac{1}{N} \sum_{i=0}^{N-1} \text{Gain} \cdot (1 - P_{total}) \tag{A5}$$

**Best Approximation/Candidate** In [48], performance profiles and the AUP aim to measure available rates across $\mathtt{m}$ tasks, where all performance metrics are threshold $\tau$, performance profiles ($\rho_{\mathtt{m}}(\tau)$) are computed as thresholds in all metrics (sorted by ascending) in current task.

$$\text{AUP}_{\mathtt{m}} = \int_{1}^{\tau_{\max}} \rho_{\mathtt{m}}(\tau) d\tau \tag{A6}$$

It is noted that the above performance profile and $\text{AUP}_{\mathtt{m}}$ score cannot measure the diversity of reasoning. Therefore, we rewrite the performance profile and AUP score:

First, the performance profile is defined in all completed trials $M_i$ in $i$-th task and the overall search space $\mathcal{H}$, as follows

$$r_i = \begin{cases} \frac{\max\{v_k : k \in \mathcal{H}\}}{v_i}, & \text{ascend} \\[2mm] \frac{v_i}{\min\{v_k : k \in \mathcal{H}\}}, & \text{descend} \end{cases} \tag{A7}$$

where ascend/descend denotes that the value is larger/smaller, the performance is better. For all trials, $\tau$ is the set of all obtained $r_i$ values.

Then, we consider all completed trials $M_i$ in $i$-th task,

$$\rho_i(\tau) = \begin{cases} |\{k \in M_i : r_k >= \tau\}|, & \text{ascend} \\ |\{k \in M_i : r_k <= \tau\}|, & \text{descend} \end{cases} \tag{A8}$$

which captures how many evaluated configurations exceed a given performance threshold $\tau$. $\rho_i(\tau)$ is the cumulative distribution curve of the trajectory.

The area term aggregates the overall concentration of strong configurations along the performance axis:

$$A = \frac{1}{N} \sum_{i=0}^{N-1} \int_{\tau_{\min}}^{\tau_{\max}} \rho_i(\tau) d\tau. \tag{A9}$$

To capture the temporal asymmetry of discovery, we compute the centroid

$$G = \frac{1}{A} \int_{\tau_{\min}}^{\tau_{\max}} x \cdot \rho_i(\tau) \, d\tau, \tag{A10}$$

and define the skewness

$$S = \int_{\tau_{\min}}^{\tau_{\max}} \left(\frac{x}{G}\right)^3 \rho(x) \, d\tau. \tag{A11}$$

Since $S$ may be unbounded and may take both positive (left-skewed) and negative (right-skewed) values, we normalize it via a reference skewness value $S_{\text{base}}$ from the GS method and a smooth monotonic mapping:

$$S' = 1 - \frac{S - S_{\text{base}}}{S_{\text{base}}}, \tag{A12}$$

$$S' = \frac{\tanh(S) + 1}{2}, S' \in (0, 1). \tag{A13}$$

Finally, the Area Under the Performance–Diversity curve $(\text{AUP}_D)$ is defined as

$$\text{AUP}_D = A/S', \tag{A14}$$

where trajectories that exhibit earlier concentration of high-performing configurations obtain larger and thus smaller $\text{AUP}_D$ value, whereas trajectories that concentrate improvements later yield smaller and therefore larger value.

## Appendix B  Details of SelfAI Benchmark

This section provides the complete technical definitions of the metrics used in the main text, including all normalization procedures, ranking aggregation, and scoring rules. These details are omitted from the main text for clarity and conciseness.

Table B1 summarizes the 12 tasks spanning six scientific categories used to evaluate the reasoning capacity of different solvers. For each dataset, we supply background

information to the LLMs to simulate the user-intent interpretation process (see Sect. A.1). Most datasets are generated through systematic grid-based exploration, whereas LCBench and the Chagas EP20 drug-discovery dataset are constructed from Bayesian optimization trajectories. LCBench, in particular, is a widely used AutoML benchmark [71]. Supplementary Fig. B2 provides an additional structural analysis of the LCBench search space. Accuracy depends on interacting hyperparameters such as weight decay, hidden-layer width, and depth, and the high-performance region is sharply localized. This structure highlights why solvers must possess strong reasoning capabilities to efficiently locate these promising configurations. Supplementary Fig. B3 reports solver rankings based on the Score metric. These rank summaries allow cross-task comparisons of solver consistency; detailed interpretation of performance trends is provided in the main text.

In Table B2, we summarize averaged performance across all benchmark tasks. Classical baselines such as grid search and TPE-based Bayesian optimization achieve competitive final objective values but receive low Score and high $AUP_D$, reflecting redundant late-stage exploration and delayed stopping. Naive LLM solvers improve early discovery but remain unstable across tasks. In contrast, SelfAI-driven solvers based on mid-sized models such as Qwen2.5-7b and GPT4-o3-mini attain the highest Scores with substantially lower $AUP_D$ and earlier stopping, indicating more focused and efficient trajectories. Larger models (e.g., Qwen2.5-72b, Llama3.3-70b) do not necessarily yield better efficiency, highlighting that structured reasoning and trajectory control within SelfAI are more critical than sheer model scale. In addition, as discussed in Sect. B.1, failure analyses are also included for completeness. These examples illustrate typical breakdowns in reasoning processes across long trajectories, such as premature stopping, limited incorporation of earlier observations, or sensitivity to minor perturbations. They inform the broader evaluation of solver stability and support the claims presented in the main text.

**Table B1**: List of tasks in SelfAI with different hyperparameters for multiple tasks and datasets.

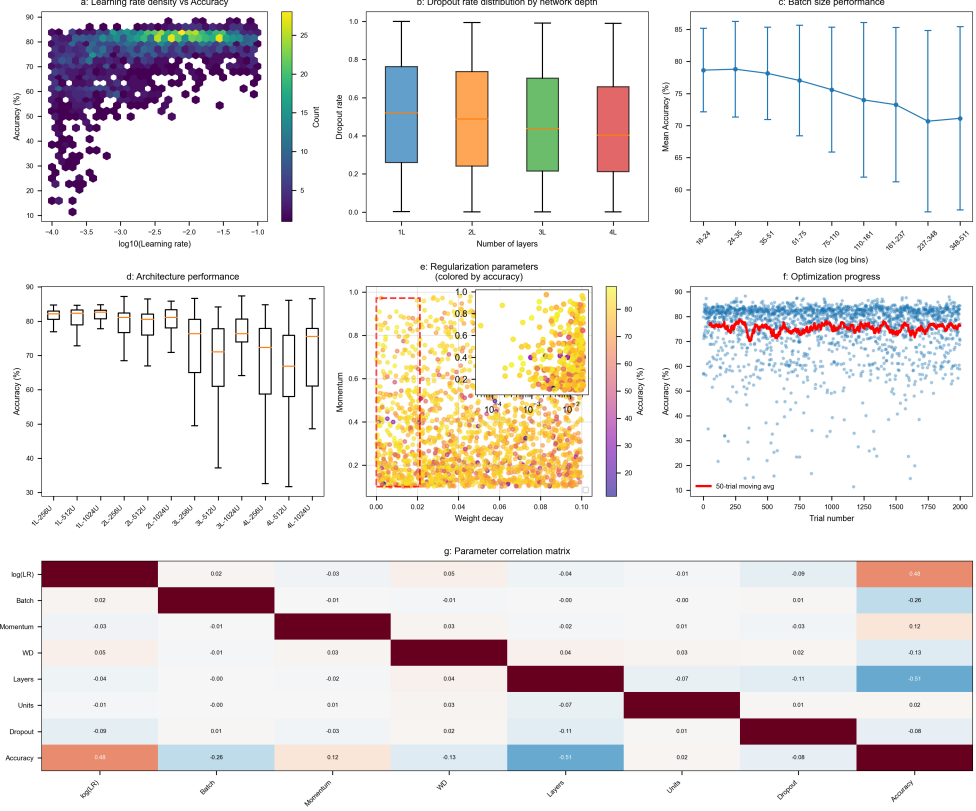| Category | Method | Task | Dim | Count | Ref. |
|---|---|---|---|---|---|
| Scientific Computing | Tensor Network | Image Completion | 3 | 64 | [65] |
| | Random Forest | Regression | 5 | 162 | [62] |
| Deep Learning | LSTM | Sentiment Analysis | 2 | 20 | [63] |
| | GraphSAGE | Node Classification | 22 | 25 | [93] |
| Computer Vision | SIREN | Image Denoising | 2 | 25 | [68] |
| | SIREN | Image Segmentation | 2 | 25 | [68] |
| | ResNet | Image Classification | 4 | 9 | [70] |
| | MAE | Image Classification | 2 | 20 | [69] |
| | FashionMnist-NN | Image Classification | 5 | 2000 | [71] |
| Medical Image Analysis | nnUnet | BraTS [82] | 3 | 18 | [72] |
| | nnUnet-revisited | BTCV [81] | 5 | 19 | [73] |
| Drug Discovery | DNN | Bioactivity Prediction | 4 | 30 | [94] |

**Fig. B2**: The relationship between different parameter settings and accuracy. **a** A hexbin plot showing the joint density of learning rate (log-scaled) and accuracy. **b** Box plots illustrating the distribution of dropout rates used for models of different depths. **c** Mean accuracy and standard error across logarithmically binned batch sizes. **d** Performance comparison of different architectural configurations (layers and units). **e** A scatter plot of weight decay against momentum, colored by accuracy. **f** The optimization trajectory, showing accuracy improvement over successive trials, with a moving average trend line. **g** A correlation matrix quantifying linear relationships between all parameters and the accuracy.

## B.1 Failure Cases

**Comparison of Best Result Hit Rate.** Table B2 compares the experimental results of traditional solvers and LLM-based solvers. Early-stopping LLM-ES partially alleviates redundant trials but remains less stable and less efficient than SelfAI-integrated variants. In contrast, SelfAI consistently enhances trajectory-level reasoning across all LLM families, reducing unnecessary exploration (lower $\text{AUP}_D$), improving stopping behavior (lower Stop-Time), and yielding higher overall Score. Mid-sized models such as Qwen2.5-7B and DeepSeek-R1-32B demonstrate particularly favorable

**Table B2**: Averaged performance comparison of SelfAI across different tasks.

| Solver | Score↑ | AUP$_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Hit-Rate↑ | Rank |
|---|---|---|---|---|---|---|---|
| GS | 0.2453 | 1.0000 | 0.5094 | 1.0000 | 1.0000 | 1.0000 | 14 |
| BS | 0.1927 | 0.8106 | 0.6265 | 0.9881 | 1.0000 | 1.0000 | 15 |
| LLM | 0.3526 | 0.7638 | 0.2949 | 1.0000 | 1.0000 | 0.9286 | 13 |
| LLM-ES | 0.5294 | 0.2349 | 0.4691 | 0.4582 | 0.9981 | 0.6429 | 3 |
| Qwen2.5-7b | 0.5562 | 0.2154 | 0.4805 | 0.3945 | 0.9957 | 0.7857 | 2 |
| Qwen2.5-14b | 0.5015 | 0.3310 | 0.4926 | 0.4997 | 0.9969 | 0.7857 | 5 |
| Qwen2.5-32b | 0.4287 | 0.4684 | 0.5252 | 0.6133 | 0.9972 | 0.7857 | 10 |
| Qwen2.5-72b | 0.4189 | 0.5358 | 0.4531 | 0.7087 | 0.9995 | 0.8571 | 11 |
| DeepSeek-r1-7b | 0.4769 | 0.1802 | 0.7020 | 0.3093 | 0.9927 | 0.5000 | 8 |
| DeepSeek-r1-14b | 0.4793 | 0.3956 | 0.4953 | 0.5100 | 0.9948 | 0.7143 | 7 |
| DeepSeek-r1-32b | 0.4989 | 0.3476 | 0.4535 | 0.5433 | 0.9933 | 0.7857 | 6 |
| DeepSeek-r1-70b | 0.4556 | 0.3513 | 0.5392 | 0.5299 | 0.9962 | 0.7143 | 9 |
| Llama3.3-70b | 0.3625 | 0.5483 | 0.5099 | 0.7271 | 0.9683 | 0.7143 | 12 |
| GPT4-o3-mini | 0.6433 | 0.2259 | 0.3168 | 0.3961 | 0.9989 | 0.8571 | 1 |
| GPT4-o3 | 0.5140 | 0.2284 | 0.5477 | 0.3961 | 0.9966 | 0.6429 | 4 |

exploration and exploitation trade-offs. Although LLMs possess advanced reasoning capabilities and can identify optimal results more efficiently, we found that LLM-based solvers still face challenges related to context limitations and computational fragility.

1) **Context limitations.** The LCBench dataset exceeds the maximum token limit of GPT-family models, which is why SelfAI does not employ GPT4-o3-mini or GPT4-o3. In practice, search spaces rarely contain thousands of trials, and this issue can be mitigated in future work. For instance, memory management and Retrieval-Augmented Generation (RAG) can help reduce token requirements by introducing additional steps to the experimental recommendation workflow, such as embedding, chunking, and retrieval.

2) **Computational fragility.** Table B2 compares the best result hit rates across solvers. The performance of LLM-based solvers is highly sensitive to the reasoning strategy, and early stopping often leads to a decline in hit rate. In the supplementary tables, we report the best result $t_{best}$ for each solver. $t_{best} = 1$ indicates that the optimal result was not found. Moreover, when $t_{best} = 1$, the score metric improves significantly, highlighting the solver's weakness in optimal stopping. As illustrated in Fig. B7, computational fragility also hinders reliable performance improvement through parameter tuning, as seen in the non-monotonic performance of the DeepSeek-R1 model series. Table C5 summarizes the metrics corresponding to Fig. B7.

Fig. B7 provides a detailed visualization of solver behavior. The Qwen2.5 models consistently identify high-performing configurations early in the search, with Qwen2.5-7B reaching near-optimal regions within the first 10–15 trials. Their trajectories exhibit stable refinement and limited oscillation, suggesting consistent integration of prior experimental evidence. In contrast, the DeepSeek-R1 models show pronounced non-monotonic behavior. DeepSeek-R1-7B frequently oscillates between high and low values, reflecting sensitivity to minor variations in its reasoning process. DeepSeek-R1-70B
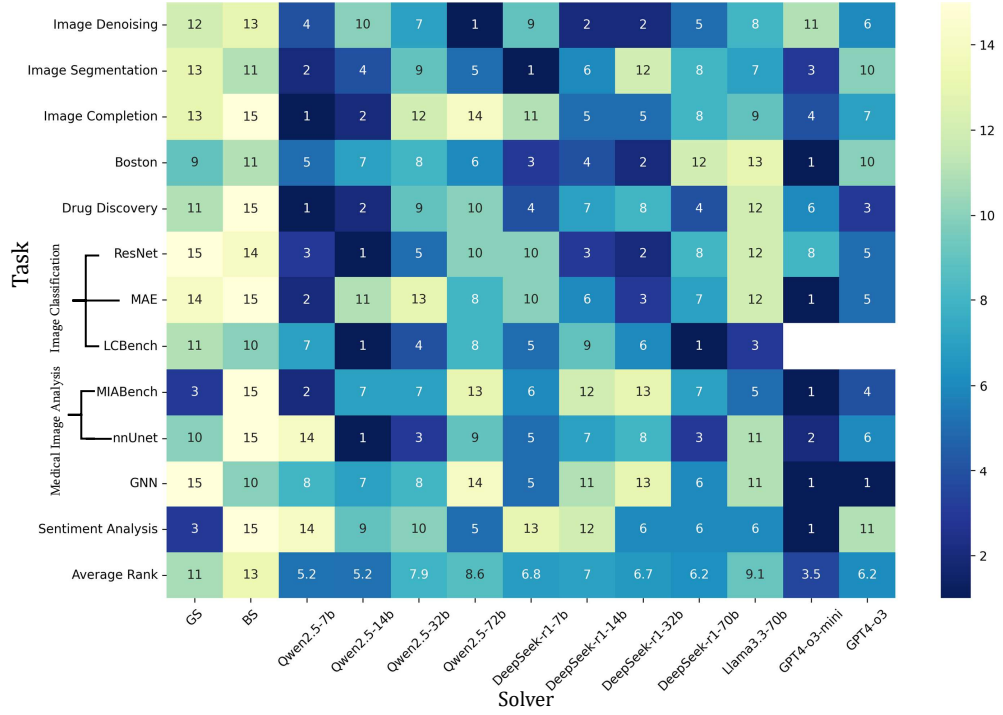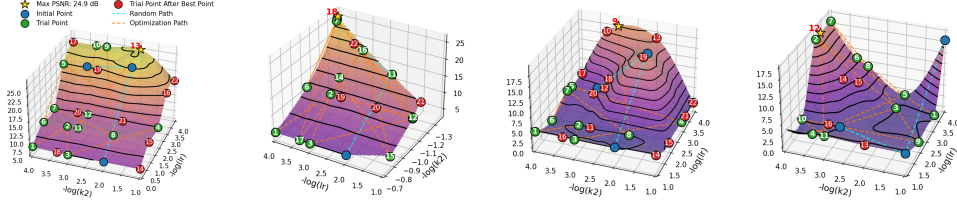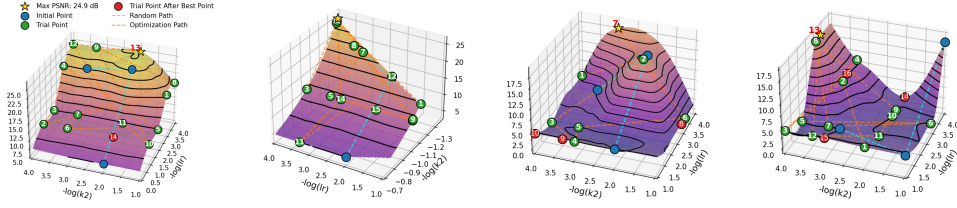
**Fig. B3**: Performance ranking of different methods across multiple tasks. The heatmap displays the rank of each method (rows) for every task (columns), where lower numbers (darker colors) indicate better performance (e.g., 1st place). The final average rank is summarized on the right.

| Task | GS | BS | Qwen2.5-7b | Qwen2.5-14b | Qwen2.5-32b | Qwen2.5-72b | DeepSeek-r1-7b | DeepSeek-r1-14b | DeepSeek-r1-32b | DeepSeek-r1-70b | Llama3.3-70b | GPT4-o3-mini | GPT4-o3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Image Denoising | 12 | 13 | 4 | 10 | 7 | 1 | 9 | 2 | 2 | 5 | 8 | 11 | 6 |
| Image Segmentation | 13 | 11 | 2 | 4 | 9 | 5 | 1 | 6 | 12 | 8 | 7 | 3 | 10 |
| Image Completion | 13 | 15 | 1 | 2 | 12 | 14 | 11 | 5 | 5 | 8 | 9 | 4 | 7 |
| Boston | 9 | 11 | 5 | 7 | 8 | 6 | 3 | 4 | 2 | 12 | 13 | 1 | 10 |
| Drug Discovery | 11 | 15 | 1 | 2 | 9 | 10 | 4 | 7 | 8 | 4 | 12 | 6 | 3 |
| ResNet | 15 | 14 | 3 | 1 | 5 | 10 | 10 | 3 | 2 | 8 | 12 | 8 | 5 |
| MAE | 14 | 15 | 2 | 11 | 13 | 8 | 10 | 6 | 3 | 7 | 12 | 1 | 5 |
| LCBench | 11 | 10 | 7 | 1 | 4 | 8 | 5 | 9 | 6 | 1 | 3 | | |
| MIABench | 3 | 15 | 2 | 7 | 7 | 13 | 6 | 12 | 13 | 7 | 5 | 1 | 4 |
| nnUnet | 10 | 15 | 14 | 1 | 3 | 9 | 5 | 7 | 8 | 3 | 11 | 2 | 6 |
| GNN | 15 | 10 | 8 | 7 | 8 | 14 | 5 | 11 | 13 | 6 | 11 | 1 | 1 |
| Sentiment Analysis | 3 | 15 | 14 | 9 | 10 | 5 | 13 | 12 | 6 | 6 | 6 | 1 | 11 |
| Average Rank | 11 | 13 | 5.2 | 5.2 | 7.9 | 8.6 | 6.8 | 7 | 6.7 | 6.2 | 9.1 | 3.5 | 6.2 |

(Image Classification spans ResNet, MAE, LCBench; Medical Image Analysis spans MIABench, nnUnet.)

achieves early improvements but then terminates prematurely, failing to further refine promising regions. Although the GS and BS reach high-performance regions, their trajectories lack the rapid breakthroughs observed in LLM-based solvers. Overall, these trajectory comparisons provide visual evidence of differences in exploration stability, breakthrough efficiency, and stopping behavior across solver families, complementing the aggregate metrics reported in the main text.
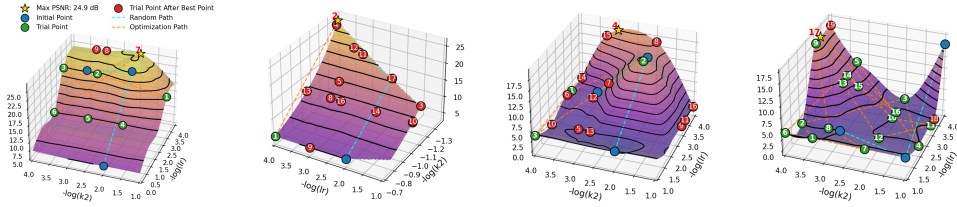
**BS**



**GPT4-o3-mini**



**GPT4-o3**



**Llama3.3-70B**



**Fig. B4**: Illustration of the optimized trajectory using the SIREN method for additional cases in image denoising (first two columns) and segmentation (last two columns). Blue points are initial points. Green points represent suggested points before reaching the optimum. Red points indicate redundant suggestions generated after the optimal point has been reached. The ⋆ marks the optimal point. The numbered labels indicate the sequence of recommendations provided by LLMs.

**Qwen2.5-7B**



**Qwen2.5-14B**



**Qwen2.5-72B**



**Qwen2.5-72B**



**Fig. B5**: More optimized trajectories of Fig. B4. These LLMs have different suggestions for non-convex hyperparameter optimization.

# DeepSeek-R1-7B



# DeepSeek-R1-14B



# DeepSeek-R1-32B



# DeepSeek-R1-70B



**Fig. B6**: More optimized trajectories of Fig. B4. These LLMs have different suggestions for non-convex hyperparameter optimization.

31

**Fig. B7**: Failure Cases for the DeepSeek-R1 Family in the scientific computing field. Both DeepSeek-R1-7B and DeepSeek-R1-70B fail to reach the optimal results, marked with a star symbol (⋆). Notably, DeepSeek-R1-7B fails to improve beyond the initial results, as indicated by the red dotted line, suggesting limited exploration capability in the search space.

# Appendix C   Supplementary Tables

**Table C3**: Performance comparison of different solvers for the Boston housing prices based on the Random Forest Regressor.

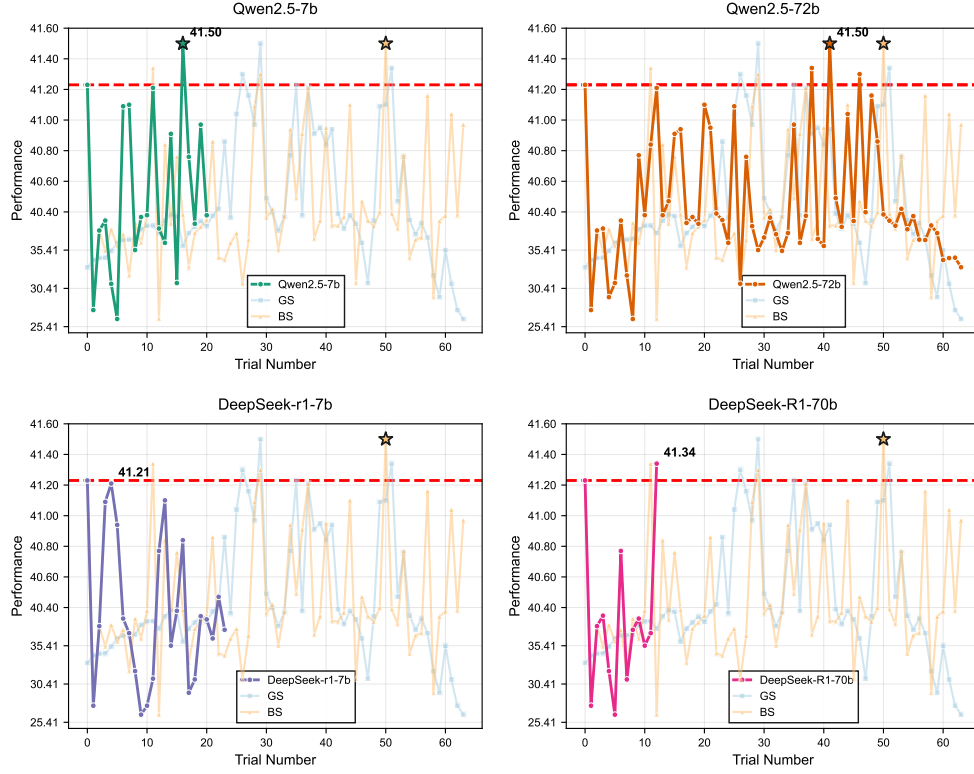| Solver | Score↑ | AUP$_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.4969 | 1.0000 | 0.0062 | 1.0000 | 0.841 | 9 |
| BS | 0.4654 | 0.9745 | 0.0755 | 0.9937 | 0.841 | 11 |
| Llama3.3-70b | 0.3831 | 0.1330 | 1.0000 | 0.1509 | 0.837 | 13 |
| Qwen2.5-7b | 0.7893 | 0.3463 | 0.0943 | 0.3270 | 0.841 | 5 |
| Qwen2.5-14b | 0.7013 | 0.6108 | 0.0063 | 0.5912 | 0.841 | 7 |
| Qwen2.5-32b | 0.6981 | 0.4081 | 0.2579 | 0.3459 | 0.841 | 8 |
| Qwen2.5-72b | 0.7862 | 0.5250 | 0.0126 | 0.4151 | 0.841 | 6 |
| DeepSeek-R1-7b | 0.9057 | 0.0910 | 0.0881 | 0.1006 | 0.841 | 3 |
| DeepSeek-R1-14b | 0.8522 | 0.3401 | 0.0063 | 0.2893 | 0.841 | 4 |
| DeepSeek-R1-32b | 0.9403 | 0.1618 | 0.0063 | 0.1132 | 0.841 | 2 |
| DeepSeek-R1-70b | 0.3948 | 0.0215 | 1.0000 | 0.0189 | 0.833 | 12 |
| GPT4-o3-mini | 0.9811 | 0.0062 | 0.0189 | 0.0189 | 0.841 | 1 |
| GPT4-o3 | 0.4694 | 0.0560 | 1.0000 | 0.0377 | 0.840 | 10 |

# Appendix D   Prompts of Compared Methods

---
**Search Prompt of LLM solver**

Instructions:
**Search Space** (Numbering starts from 0, excluding Completed Trials):
{{TRIALS}}
Task 1: Optimization Recommendation Recommend exactly {{N_JOBS}} promising trials from the provided **Search Space** (include both number and params).
**Rules:**
1. "params" MUST include:
 {{HYPERNAME}}
2. All selected 'params' must match exactly with the provided **Search Space**. Do NOT leave out any key.
3. Use the analysis in **Task 1** (trial analysis, performance trends, highlights, and other insights) to guide selection.

---

**Table C4**: Performance comparison of different solvers for the sentiment analysis task based on the LSTM model.

| Solver | Score↑ | AUP$_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.4750 | 1.0000 | 0.0500 | 1.0000 | 0.96 | 2 |
| BS | 0.0294 | 0.8817 | 0.9412 | 1.0000 | 0.96 | 13 |
| Llama3.3-70b | 0.3824 | 0.8817 | 0.2353 | 1.0000 | 0.96 | 4 |
| Qwen2.5-7b | 0.1830 | 0.2857 | 1.0000 | 0.5294 | 0.94 | 12 |
| Qwen2.5-14b | 0.3529 | 0.8318 | 0.4118 | 0.8824 | 0.96 | 7 |
| Qwen2.5-32b | 0.3529 | 0.6383 | 0.5882 | 0.7059 | 0.96 | 8 |
| Qwen2.5-72b | 0.4412 | 0.8220 | 0.2353 | 0.8824 | 0.96 | 3 |
| DeepSeek-R1-7b | 0.2157 | 0.1419 | 1.0000 | 0.3529 | 0.93 | 11 |
| DeepSeek-R1-14b | 0.2353 | 0.8817 | 0.5294 | 1.0000 | 0.96 | 10 |
| DeepSeek-R1-32b | 0.3824 | 0.8817 | 0.2353 | 1.0000 | 0.96 | 4 |
| DeepSeek-R1-70b | 0.3824 | 0.8817 | 0.2353 | 1.0000 | 0.96 | 4 |
| GPT4-o3-mini | 0.8824 | 0.1494 | 0.0588 | 0.1765 | 0.96 | 1 |
| GPT4-o3 | 0.2745 | 0.0469 | 1.0000 | 0.1765 | 0.93 | 9 |

4. Based on the above analysis, explore under-explored regions only when there is clear evidence of potential performance gain.
5. Do not mix, modify, or create new values.
6. You MUST not output any JSON blocks in this part.
7. You MUST provide reasoning for each recommendation.

**Early-Stopping Prompt of LLM-ES solver**

Completed trials: `{{COMPLETED_TRIALS}}`
The following **Search Space** contains **unexplored** trials. `{{TRIALS}}`
If the optimization process should be stopped, Answer: Yes with confidence score: `{{CONFIDENCE_SOCRE}}`. Otherwise, Answer: No with confidence score: `{{CONFIDENCE_SOCRE}}`.
Finally, you MUST output 'Answer: No/Yes' with confidence score: `{{CONFIDENCE_SOCRE}}`.

# References

[1] OpenAI, R.: Gpt-4 technical report. arxiv 2303.08774. View in Article **2**(5) (2023)

**Table C5**: Performance comparison of different solvers for scientific computing fields, where we evaluate the tensor wheel decomposition method on the multi-spectral image completion dataset.

| Solver | Score↑ | $\text{AUP}_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.2656 | 1.0000 | 0.4688 | 1.0000 | 41.50 | 11 |
| BS | 0.1066 | 0.9638 | 0.7869 | 1.0000 | 41.50 | 13 |
| Llama3.3-70b | 0.4016 | 0.8979 | 0.2623 | 0.9344 | 41.50 | 8 |
| Qwen2.5-7b | 0.7377 | 0.2904 | 0.2295 | 0.2951 | 41.50 | 1 |
| Qwen2.5-14b | 0.5410 | 0.2779 | 0.4262 | 0.4918 | 41.50 | 2 |
| Qwen2.5-32b | 0.2869 | 0.9638 | 0.4262 | 1.0000 | 41.50 | 10 |
| Qwen2.5-72b | 0.1803 | 0.9638 | 0.6393 | 1.0000 | 41.50 | 12 |
| DeepSeek-R1-7b | 0.3216 | 0.3009 | 1.0000 | 0.3443 | 41.21 | 9 |
| DeepSeek-R1-14b | 0.4426 | 0.8687 | 0.2295 | 0.8852 | 41.50 | 4 |
| DeepSeek-R1-32b | 0.4426 | 0.6645 | 0.4426 | 0.6721 | 41.50 | 4 |
| DeepSeek-R1-70b | 0.4136 | 0.1477 | 1.0000 | 0.1639 | 41.34 | 7 |
| GPT4-o3-mini | 0.5410 | 0.2623 | 0.4590 | 0.4590 | 41.50 | 3 |
| GPT4-o3 | 0.4217 | 0.0539 | 1.0000 | 0.1475 | 41.34 | 6 |

[2] Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., et al.: Qwen technical report. arXiv preprint arXiv:2309.16609 (2023)

[3] Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al.: Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. arXiv preprint arXiv:2501.12948 (2025)

[4] Huang, J., Chang, K.C.-C.: Towards reasoning in large language models: A survey. arXiv preprint arXiv:2212.10403 (2022)

[5] Kojima, T., Gu, S.S., Reid, M., Matsuo, Y., Iwasawa, Y.: Large language models are zero-shot reasoners. Advances in neural information processing systems **35**, 22199–22213 (2022)

[6] Wu, S., Fei, H., Qu, L., Ji, W., Chua, T.-S.: Next-gpt: Any-to-any multimodal llm. In: Forty-first International Conference on Machine Learning (2024)

[7] McKinzie, B., Gan, Z., Fauconnier, J.-P., Dodge, S., Zhang, B., Dufter, P., Shah, D., Du, X., Peng, F., Belyi, A., *et al.*: Mm1: methods, analysis and insights from multimodal llm pre-training. In: European Conference on Computer Vision, pp. 304–323 (2024). Springer

[8] Ferrag, M.A., Tihanyi, N., Debbah, M.: From llm reasoning to autonomous ai

**Table C6**: Performance comparison of different solvers for the image segmentation task based on the SIREN model.

| Solver | Score↑ | AUP$_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.1100 | 1.000000 | 0.780000 | 1.000000 | 16.63 | 13 |
| BS | 0.329545 | 0.761200 | 0.477273 | 0.863636 | 16.63 | 11 |
| Qwen2.5-7b | 0.602273 | 0.487619 | 0.363636 | 0.431818 | 16.63 | 2 |
| Qwen2.5-14b | 0.454545 | 0.708497 | 0.272727 | 0.818182 | 16.63 | 4 |
| Qwen2.5-32b | 0.374388 | 0.538240 | 0.613636 | 0.636364 | 16.60 | 9 |
| Qwen2.5-72b | 0.408757 | 0.536178 | 0.636364 | 0.545455 | 16.60 | 5 |
| DeepSeek-r1-7b | 0.693182 | 0.318656 | 0.227273 | 0.386364 | 16.63 | 1 |
| DeepSeek-r1-14b | 0.408534 | 0.471791 | 0.704545 | 0.477273 | 16.60 | 6 |
| DeepSeek-r1-32b | 0.312477 | 0.198417 | 1.000000 | 0.340909 | 15.77 | 12 |
| DeepSeek-r1-70b | 0.374555 | 0.576330 | 0.590909 | 0.659091 | 16.60 | 8 |
| Llama3.3-70b | 0.380344 | 0.285865 | 0.590909 | 0.295455 | 10.79 | 7 |
| GPT4-o3-mini | 0.477273 | 0.517709 | 0.454545 | 0.590909 | 16.63 | 3 |
| GPT4-o3 | 0.363636 | 0.671366 | 0.477273 | 0.795455 | 16.63 | 10 |

agents: A comprehensive review. arXiv preprint arXiv:2504.19678 (2025)

[9] Ghafarollahi, A., Buehler, M.J.: Sciagents: automating scientific discovery through bioinspired multi-agent intelligent graph reasoning. Advanced Materials **37**(22), 2413523 (2025)

[10] Chen, D., Bai, Y., Ament, S., Zhao, W., Guevarra, D., Zhou, L., Selman, B., Dover, R.B., Gregoire, J.M., Gomes, C.P.: Automating crystal-structure phase mapping by combining deep learning with constraint reasoning. Nature Machine Intelligence **3**(9), 812–822 (2021)

[11] Lin, J., Guo, Y., Han, Y., Hu, S., Ni, Z., Wang, L., Chen, M., Liu, H., Chen, R., He, Y., et al.: Se-agent: Self-evolution trajectory optimization in multi-step reasoning with llm-based agents. arXiv preprint arXiv:2508.02085 (2025)

[12] Qin, Y., Liang, S., Ye, Y., Zhu, K., Yan, L., Lu, Y., Lin, Y., Cong, X., Tang, X., Qian, B., et al.: Toolllm: Facilitating large language models to master 16000+ real-world apis. arXiv preprint arXiv:2307.16789 (2023)

[13] Xu, J., Du, W., Liu, X., Li, X.: Llm4workflow: An llm-based automated workflow model generation tool. In: Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering, pp. 2394–2398 (2024)

**Table C7**: Performance comparison of different solvers for the image denoising task based on the SIREN model.

| Solver | Score↑ | AUP$_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.1500 | 1.0000 | 0.7000 | 1.0000 | 24.78 | 12 |
| BS | 0.1477 | 0.8819 | 0.7045 | 1.0000 | 24.78 | 13 |
| Qwen2.5-7b | 0.7045 | 0.1802 | 0.1364 | 0.4545 | 24.78 | 4 |
| Llama3.3-70b | 0.5706 | 0.1105 | 0.5455 | 0.2955 | 24.27 | 8 |
| Qwen2.5-14b | 0.4205 | 0.5205 | 0.5455 | 0.6136 | 24.78 | 10 |
| Qwen2.5-32b | 0.5795 | 0.4762 | 0.3182 | 0.5227 | 24.78 | 7 |
| Qwen2.5-72b | 0.7614 | 0.1543 | 0.1364 | 0.3409 | 24.78 | 1 |
| DeepSeek-R1-7b | 0.5341 | 0.4193 | 0.3864 | 0.5455 | 24.78 | 9 |
| DeepSeek-R1-14b | 0.7159 | 0.1639 | 0.0909 | 0.4773 | 24.78 | 2 |
| DeepSeek-R1-32b | 0.7159 | 0.1284 | 0.2045 | 0.3636 | 24.78 | 2 |
| DeepSeek-R1-70b | 0.6818 | 0.1359 | 0.2955 | 0.3409 | 24.78 | 5 |
| GPT4-o3-mini | 0.2703 | 0.4389 | 0.7955 | 0.6591 | 24.48 | 11 |
| GPT4-o3 | 0.6023 | 0.4786 | 0.2045 | 0.5909 | 24.78 | 6 |

[14] Gupta, S., Mahmood, A., Shetty, P., Adeboye, A., Ramprasad, R.: Data extraction from polymer literature using large language models. Communications materials **5**(1), 269 (2024)

[15] Wang, H., Fu, T., Du, Y., Gao, W., Huang, K., Liu, Z., Chandak, P., Liu, S., Van Katwyk, P., Deac, A., *et al.*: Scientific discovery in the age of artificial intelligence. Nature **620**(7972), 47–60 (2023)

[16] Darvish, K., Skreta, M., Zhao, Y., Yoshikawa, N., Som, S., Bogdanovic, M., Cao, Y., Hao, H., Xu, H., Aspuru-Guzik, A., et al.: Organa: A robotic assistant for automated chemistry experimentation and characterization. Matter **8**(2) (2025)

[17] Huang, K., Zhang, S., Wang, H., Qu, Y., Lu, Y., Roohani, Y., Li, R., Qiu, L., Zhang, J., Di, Y., et al.: Biomni: A general-purpose biomedical ai agent. bioRxiv, 2025–05 (2025)

[18] Alampara, N., Schilling-Wilhelmi, M., Ríos-García, M., Mandal, I., Khetarpal, P., Grover, H.S., Krishnan, N.A., Jablonka, K.M.: Probing the limitations of multimodal language models for chemistry and materials research. Nature computational science, 1–10 (2025)

[19] Polak, M.P., Morgan, D.: Extracting accurate materials data from research papers with conversational language models and prompt engineering. Nature Communications **15**(1), 1569 (2024)

**Table C8**: Performance comparison of different solvers for Mask AutoEncoders (MAE) [69].

| Solver | Score↑ | $AUP_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.1000 | 1.0000 | 0.8000 | 1.0000 | 85.0 | 12 |
| BS | 0.0294 | 0.8093 | 0.9412 | 1.0000 | 85.0 | 13 |
| Llama3.3-70b | 0.2941 | 0.8093 | 0.4118 | 1.0000 | 85.0 | 10 |
| Qwen2.5-7b | 0.7647 | 0.0112 | 0.2353 | 0.2353 | 85.0 | 2 |
| Qwen2.5-14b | 0.3471 | 0.1513 | 1.0000 | 0.2941 | 84.5 | 9 |
| Qwen2.5-32b | 0.2647 | 0.8093 | 0.4706 | 1.0000 | 85.0 | 11 |
| Qwen2.5-72b | 0.4412 | 0.8093 | 0.1176 | 1.0000 | 85.0 | 7 |
| DeepSeek-R1-7b | 0.4104 | 0.0080 | 1.0000 | 0.1765 | 84.9 | 8 |
| DeepSeek-R1-14b | 0.6176 | 0.3853 | 0.2941 | 0.4706 | 85.0 | 5 |
| DeepSeek-R1-32b | 0.6765 | 0.1108 | 0.2941 | 0.3529 | 85.0 | 3 |
| DeepSeek-R1-70b | 0.5588 | 0.2235 | 0.1765 | 0.7059 | 85.0 | 6 |
| GPT4-o3-mini | 0.8235 | 0.0088 | 0.1176 | 0.2353 | 85.0 | 1 |
| GPT4-o3 | 0.6471 | 0.1141 | 0.3529 | 0.3529 | 85.0 | 4 |

[20] Dagdelen, J., Dunn, A., Lee, S., Walker, N., Rosen, A.S., Ceder, G., Persson, K.A., Jain, A.: Structured information extraction from scientific text with large language models. Nature communications **15**(1), 1418 (2024)

[21] Mu, C., Zhang, X., Wang, H.: Planning of heuristics: Strategic planning on large language models with monte carlo tree search for automating heuristic optimization. arXiv preprint arXiv:2502.11422 (2025)

[22] Zhou, A., Yan, K., Shlapentokh-Rothman, M., Wang, H., Wang, Y.-X.: Language agent tree search unifies reasoning acting and planning in language models. arXiv preprint arXiv:2310.04406 (2023)

[23] Hao, S., Gu, Y., Ma, H., Hong, J.J., Wang, Z., Wang, D.Z., Hu, Z.: Reasoning with language model is planning with world model. arXiv preprint arXiv:2305.14992 (2023)

[24] Toledo, E., Hambardzumyan, K., Josifoski, M., Hazra, R., Baldwin, N., Audran-Reiss, A., Kuchnik, M., Magka, D., Jiang, M., Lupidi, A.M., et al.: Ai research agents for machine learning: Search, exploration, and generalization in mle-bench. arXiv preprint arXiv:2507.02554 (2025)

[25] Zweiger, A., Pari, J., Guo, H., Akyürek, E., Kim, Y., Agrawal, P.: Self-adapting language models. arXiv preprint arXiv:2506.10943 (2025)

**Table C9**: Performance comparison of different solvers for the image classification task on the ImageNet dataset reported from the ResNet benchmark.

| Solver | Score↑ | AUP$_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.0000 | 1.0000 | 1.0000 | 1.0000 | 21.43 | |
| BS | 0.0833 | 0.6394 | 0.8333 | 1.0000 | 21.43 | |
| Llama3.3-70b | 0.3333 | 0.6394 | 0.3333 | 1.0000 | 21.43 | |
| Qwen2.5-7b | 0.6667 | 0.0865 | 0.3333 | 0.3333 | 21.43 | |
| Qwen2.5-14b | 0.8333 | 0.0000 | 0.1667 | 0.1667 | 21.43 | |
| Qwen2.5-32b | 0.5833 | 0.2192 | 0.3333 | 0.5000 | 21.43 | |
| Qwen2.5-72b | 0.4167 | 0.2988 | 0.5000 | 0.6667 | 21.43 | |
| DeepSeek-R1-7b | 0.4167 | 0.3322 | 0.5000 | 0.6667 | 21.43 | |
| DeepSeek-R1-14b | 0.6667 | 0.0865 | 0.3333 | 0.3333 | 21.43 | |
| DeepSeek-R1-32b | 0.7500 | 0.0865 | 0.1667 | 0.3333 | 21.43 | |
| DeepSeek-R1-70b | 0.5000 | 0.2443 | 0.5000 | 0.5000 | 21.43 | |
| GPT4-o3-mini | 0.5000 | 0.2192 | 0.5000 | 0.5000 | 21.43 | |
| GPT4-o3 | 0.5833 | 0.2192 | 0.3333 | 0.5000 | 21.43 | |

[26] Team, N., Zhang, B., Feng, S., Yan, X., Yuan, J., Yu, Z., He, X., Huang, S., Hou, S., Nie, Z., et al.: Novelseek: When agent becomes the scientist–building closed-loop system from hypothesis to verification. arXiv preprint arXiv:2505.16938 (2025)

[27] Baek, J., Jauhar, S.K., Cucerzan, S., Hwang, S.J.: Researchagent: Iterative research idea generation over scientific literature with large language models. In: Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pp. 6709–6738 (2025)

[28] Steyaert, S., Pizurica, M., Nagaraj, D., Khandelwal, P., Hernandez-Boussard, T., Gentles, A.J., Gevaert, O.: Multimodal data fusion for cancer biomarker discovery with deep learning. Nature machine intelligence **5**(4), 351–362 (2023)

[29] Gao, F., Li, H., Chen, Z., Yi, Y., Nie, S., Cheng, Z., Liu, Z., Guo, Y., Liu, S., Qin, Q., *et al.*: A chemical autonomous robotic platform for end-to-end synthesis of nanoparticles. Nature Communications **16**(1), 7558 (2025)

[30] Zhang, Y., Han, Y., Chen, S., Yu, R., Zhao, X., Liu, X., Zeng, K., Yu, M., Tian, J., Zhu, F., et al.: Large language models to accelerate organic chemistry synthesis. Nature Machine Intelligence, 1–13 (2025)

**Table C10**: Performance comparison of different solvers for the image classification task on the LCBench dataset reported from [71].

| Solver | Score↑ | AUP$_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.1750 | 1.0000 | 0.6500 | 1.0000 | 88.29 | |
| BS | 0.4189 | 0.9985 | 0.1622 | 1.0000 | 88.29 | |
| Llama3.3-70b | 0.4738 | 0.0052 | 1.0000 | 0.0486 | 87.98 | |
| Qwen2.5-7b | 0.4619 | 0.0047 | 1.0000 | 0.0451 | 85.78 | |
| Qwen2.5-14b | 0.4739 | 0.0017 | 1.0000 | 0.0180 | 85.62 | |
| Qwen2.5-32b | 0.4718 | 0.0022 | 1.0000 | 0.0225 | 85.62 | |
| Qwen2.5-72b | 0.4599 | 0.0081 | 1.0000 | 0.0766 | 87.98 | |
| DeepSeek-R1-7b | 0.4691 | 0.0027 | 1.0000 | 0.0280 | 85.62 | |
| DeepSeek-R1-14b | 0.4579 | 0.0086 | 1.0000 | 0.0806 | 87.98 | |
| DeepSeek-R1-32b | 0.4677 | 0.0035 | 1.0000 | 0.0331 | 85.78 | |
| DeepSeek-R1-70b | 0.4739 | 0.0014 | 1.0000 | 0.0180 | 85.62 | |
| GPT4-o3-mini | - | - | - | - | - | - |
| GPT4-o3 | - | - | - | - | - | - |

[31] Xiong, J., Zhang, W., Wang, Y., Huang, J., Shi, Y., Xu, M., Li, M., Fu, Z., Kong, X., Wang, Y., et al.: Bridging chemistry and artificial intelligence by a reaction description language. Nature Machine Intelligence, 1–12 (2025)

[32] Jablonka, K.M., Schwaller, P., Ortega-Guerrero, A., Smit, B.: Leveraging large language models for predictive chemistry. Nature Machine Intelligence **6**(2), 161–169 (2024)

[33] Miret, S., Krishnan, N.A.: Enabling large language models for real-world materials discovery. Nature Machine Intelligence, 1–8 (2025)

[34] Kang, Y., Kim, J.: Chatmof: an artificial intelligence system for predicting and generating metal-organic frameworks using large language models. Nature communications **15**(1), 4705 (2024)

[35] Stach, E., DeCost, B., Kusne, A.G., Hattrick-Simpers, J., Brown, K.A., Reyes, K.G., Schrier, J., Billinge, S., Buonassisi, T., Foster, I., *et al.*: Autonomous experimentation systems for materials development: A community perspective. Matter **4**(9), 2702–2726 (2021)

[36] Boiko, D.A., MacKnight, R., Kline, B., Gomes, G.: Autonomous chemical research with large language models. Nature **624**(7992), 570–578 (2023)

[37] Özçelik, R., Ruiter, S., Criscuolo, E., Grisoni, F.: Chemical language modeling

**Table C11**: Performance comparison of different solvers for the medical image segmentation task on the BTCV dataset reported from the nnUnet benchmark [73].

| Solver | Score↑ | AUP$_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.4211 | 1.0000 | 0.1579 | 1.0000 | 85.04 | 3 |
| BS | 0.0625 | 0.6453 | 0.8750 | 1.0000 | 85.04 | 13 |
| Llama3.3-70b | 0.3750 | 0.6453 | 0.2500 | 1.0000 | 85.04 | 5 |
| Qwen2.5-7b | 0.4688 | 0.3441 | 0.5000 | 0.5625 | 85.04 | 2 |
| Qwen2.5-14b | 0.2813 | 0.6453 | 0.4375 | 1.0000 | 85.04 | 7 |
| Qwen2.5-32b | 0.2813 | 0.4742 | 0.6875 | 0.7500 | 85.04 | 7 |
| Qwen2.5-72b | 0.1875 | 0.6453 | 0.6250 | 1.0000 | 85.04 | 11 |
| DeepSeek-R1-7b | 0.3680 | 0.0813 | 1.0000 | 0.1250 | 83.69 | 6 |
| DeepSeek-R1-14b | 0.1983 | 0.0505 | 1.0000 | 0.1875 | 80.69 | 10 |
| DeepSeek-R1-32b | 0.1875 | 0.6453 | 0.6250 | 1.0000 | 85.04 | 11 |
| DeepSeek-R1-70b | 0.2813 | 0.6453 | 0.4375 | 1.0000 | 85.04 | 7 |
| GPT4-o3-mini | 0.6875 | 0.1466 | 0.2500 | 0.3750 | 85.04 | 1 |
| GPT4-o3 | 0.4063 | 0.5081 | 0.4375 | 0.7500 | 85.04 | 4 |

with structured state space sequence models. Nature Communications **15**(1), 6176 (2024)

[38] Ghareeb, A.E., Chang, B., Mitchener, L., Yiu, A., Szostkiewicz, C.J., Laurent, J.M., Razzak, M.T., White, A.D., Hinks, M.M., Rodriques, S.G.: Robin: A multi-agent system for automating scientific discovery. arXiv preprint arXiv:2505.13400 (2025)

[39] Abramson, J., Adler, J., Dunger, J., Evans, R., Green, T., Pritzel, A., Ronneberger, O., Willmore, L., Ballard, A.J., Bambrick, J., *et al.*: Accurate structure prediction of biomolecular interactions with alphafold 3. Nature **630**(8016), 493–500 (2024)

[40] Brahmavar, S.B., Srinivasan, A., Dash, T., Krishnan, S.R., Vig, L., Roy, A., Aduri, R.: Generating novel leads for drug discovery using llms with logical feedback. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol. 38, pp. 21–29 (2024)

[41] Jiménez-Luna, J., Grisoni, F., Schneider, G.: Drug discovery with explainable artificial intelligence. Nature Machine Intelligence **2**(10), 573–584 (2020)

[42] King, R.D., Whelan, K.E., Jones, F.M., Reiser, P.G., Bryant, C.H., Muggleton, S.H., Kell, D.B., Oliver, S.G.: Functional genomic hypothesis generation and

**Table C12**: Performance comparison of different solvers for the nnUnet model on the BraTS dataset [82].

| Solver | Score↑ | AUP$_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.1667 | 1.0000 | 0.6667 | 1.0000 | 82.45 | 10 |
| BS | 0.0333 | 0.3786 | 0.9333 | 1.0000 | 82.45 | 13 |
| Llama3.3-70b | 0.1333 | 0.3786 | 0.7333 | 1.0000 | 82.45 | 11 |
| Qwen2.5-7b | 0.1000 | 0.3786 | 0.8000 | 1.0000 | 82.45 | 12 |
| Qwen2.5-14b | 0.4333 | 0.0698 | 0.5333 | 0.6000 | 82.45 | 1 |
| Qwen2.5-32b | 0.3333 | 0.3118 | 0.5333 | 0.8000 | 82.45 | 3 |
| Qwen2.5-72b | 0.2333 | 0.0883 | 0.7333 | 0.8000 | 82.45 | 9 |
| DeepSeek-R1-7b | 0.2995 | 0.0481 | 1.0000 | 0.4000 | 82.41 | 5 |
| DeepSeek-R1-14b | 0.2940 | 0.0406 | 1.0000 | 0.4000 | 82.00 | 7 |
| DeepSeek-R1-32b | 0.2667 | 0.0839 | 0.6667 | 0.8000 | 82.45 | 8 |
| DeepSeek-R1-70b | 0.3333 | 0.3126 | 0.5333 | 0.8000 | 82.45 | 3 |
| GPT4-o3-mini | 0.4000 | 0.2347 | 0.6000 | 0.6000 | 82.45 | 2 |
| GPT4-o3 | 0.2984 | 0.0484 | 1.0000 | 0.4000 | 82.33 | 6 |

experimentation by a robot scientist. Nature **427**(6971), 247–252 (2004)

[43] Mandal, I., Soni, J., Zaki, M., Smedskjaer, M.M., Wondraczek, K., Wondraczek, L., Gosvami, N.N., Krishnan, N.A.: Evaluating large language model agents for automation of atomic force microscopy. Nature Communications **16**(1), 9104 (2025)

[44] Audran-Reiss, A., EstapÃŠ, J.A., Hambardzumyan, K., Budhiraja, A., Josifoski, M., Toledo, E., Hazra, R., Magka, D., Shvartsman, M., Pathak, P., et al.: What does it take to be a good ai research agent? studying the role of ideation diversity. arXiv preprint arXiv:2511.15593 (2025)

[45] Gottweis, J., Weng, W.-H., Daryin, A., Tu, T., Palepu, A., Sirkovic, P., Myaskovsky, A., Weissenberger, F., Rong, K., Tanno, R., et al.: Towards an ai co-scientist. arXiv preprint arXiv:2502.18864 (2025)

[46] Li, X., Wang, S., Zeng, S., Wu, Y., Yang, Y.: A survey on llm-based multi-agent systems: workflow, infrastructure, and challenges. Vicinagearth **1**(1), 9 (2024)

[47] Yamada, Y., Lange, R.T., Lu, C., Hu, S., Lu, C., Foerster, J., Clune, J., Ha, D.: The ai scientist-v2: Workshop-level automated scientific discovery via agentic tree search. arXiv preprint arXiv:2504.08066 (2025)

[48] Nathani, D., Madaan, L., Roberts, N., Bashlykov, N., Menon, A., Moens, V.,

**Table C13**: Performance comparison of different solvers for the Graph-SAGE model on the imbalanced node classification task.

| Solver | Score↑ | AUP$_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.1000 | 1.0000 | 0.8000 | 1.0000 | 89.28 | 13 |
| BS | 0.2727 | 0.8982 | 0.4545 | 1.0000 | 89.28 | 8 |
| Llama3.3-70b | 0.1591 | 0.8982 | 0.6818 | 1.0000 | 89.28 | 9 |
| Qwen2.5-7b | 0.2880 | 0.0928 | 1.0000 | 0.4091 | 89.12 | 6 |
| Qwen2.5-14b | 0.2894 | 0.0890 | 1.0000 | 0.4091 | 89.15 | 5 |
| Qwen2.5-32b | 0.2880 | 0.3564 | 1.0000 | 0.4091 | 89.12 | 6 |
| Qwen2.5-72b | 0.1364 | 0.8982 | 0.7273 | 1.0000 | 89.28 | 12 |
| DeepSeek-R1-7b | 0.3116 | 0.3083 | 1.0000 | 0.3636 | 89.15 | 3 |
| DeepSeek-R1-14b | 0.1591 | 0.8982 | 0.6818 | 1.0000 | 89.28 | 9 |
| DeepSeek-R1-32b | 0.1591 | 0.8657 | 0.7273 | 0.9545 | 89.28 | 11 |
| DeepSeek-R1-70b | 0.2955 | 0.7096 | 0.5909 | 0.8182 | 89.28 | 4 |
| GPT4-o3-mini | 0.7727 | 0.1463 | 0.0455 | 0.4091 | 89.282 | 1 |
| GPT4-o3 | 0.7727 | 0.2540 | 0.1818 | 0.2727 | 89.282 | 1 |

Budhiraja, A., Magka, D., Vorotilov, V., Chaurasia, G., et al.: Mlgym: A new framework and benchmark for advancing ai research agents. arXiv preprint arXiv:2502.14499 (2025)

[49] Huang, Q., Vora, J., Liang, P., Leskovec, J.: Mlagentbench: Evaluating language agents on machine learning experimentation. arXiv preprint arXiv:2310.03302 (2023)

[50] Liao, R., Qiu, J., Chen, X., Li, X.: Llm4eo: Large language model for evolutionary optimization in flexible job shop scheduling. arXiv preprint arXiv:2511.16485 (2025)

[51] Jiang, Q., Karniadakis, G.: Agenticsciml: Collaborative multi-agent systems for emergent discovery in scientific machine learning. arXiv preprint arXiv:2511.07262 (2025)

[52] Meng, S., Wang, Y., Yang, C.-F., Peng, N., Chang, K.-W.: Llm-a*: Large language model enhanced incremental heuristic search on path planning. In: EMNLP (Findings) (2024)

[53] Kochnev, R., Goodarzi, A.T., Bentyn, Z.A., Ignatov, D., Timofte, R.: Optuna vs code llama: Are llms a new paradigm for hyperparameter tuning? arXiv preprint arXiv:2504.06006 (2025)

**Table C14**: Performance comparison of different solvers for bioactivity prediction on Chagas EP20 dataset [95]

| Solver | Score↑ | AUP$_D$ ↓ | Best-Time↓ | Stop-Time↓ | Best Result↑ | Rank |
|---|---|---|---|---|---|---|
| GS | 0.4833 | 1.0000 | 0.0333 | 1.0000 | 0.754 | 11 |
| BS | 0.3333 | 0.8941 | 0.3333 | 1.0000 | 0.754 | 13 |
| Llama3.3-70b | 0.4630 | 0.8941 | 0.0741 | 1.0000 | 0.754 | 12 |
| Qwen2.5-7b | 0.9074 | 0.0772 | 0.0741 | 0.1111 | 0.754 | 1 |
| Qwen2.5-14b | 0.8889 | 0.0649 | 0.1111 | 0.1111 | 0.754 | 2 |
| Qwen2.5-32b | 0.6296 | 0.4231 | 0.0741 | 0.6667 | 0.754 | 9 |
| Qwen2.5-72b | 0.5741 | 0.6799 | 0.0741 | 0.7778 | 0.754 | 10 |
| DeepSeek-R1-7b | 0.7778 | 0.1104 | 0.2222 | 0.2222 | 0.754 | 4 |
| DeepSeek-R1-14b | 0.7037 | 0.5510 | 0.0741 | 0.5185 | 0.754 | 7 |
| DeepSeek-R1-32b | 0.6852 | 0.3405 | 0.0741 | 0.5556 | 0.754 | 8 |
| DeepSeek-R1-70b | 0.7778 | 0.3160 | 0.1111 | 0.3333 | 0.754 | 4 |
| GPT4-o3-mini | 0.7407 | 0.3546 | 0.1852 | 0.3333 | 0.754 | 6 |
| GPT4-o3 | 0.8148 | 0.0622 | 0.0370 | 0.3333 | 0.754 | 3 |

[54] Hao, S., Gu, Y., Ma, H., Hong, J.J., Wang, Z., Wang, D.Z., Hu, Z.: Reasoning with language model is planning with world model. arXiv preprint arXiv:2305.14992 (2023)

[55] Park, J.S., O'Brien, J.C., Cai, C.J., Morris, M.R., Liang, P.: Generative agents: Interactive simulacra of human behavior. arXiv preprint arXiv:2304.03442 (2023)

[56] Akiba, T., Sano, S., Yanase, T., Ohta, T., Koyama, M.: Optuna: A next-generation hyperparameter optimization framework. In: The 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp. 2623–2631 (2019)

[57] Yadan, O.: Hydra - A framework for elegantly configuring complex applications. Github (2019). https://github.com/facebookresearch/hydra

[58] Feurer, M., Klein, A., Eggensperger, K., Springenberg, J., Blum, M., Hutter, F.: Efficient and robust automated machine learning. In: Advances in Neural Information Processing Systems 28 (2015), pp. 2962–2970 (2015)

[59] Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023)

[60] Watanabe, S.: Tree-structured parzen estimator: Understanding its algorithm

components and their roles for better empirical performance. arXiv preprint arXiv:2304.11127 (2023)

[61] Akiba, T., Sano, S., Yanase, T., Ohta, T., Koyama, M.: Optuna: A next-generation hyperparameter optimization framework. In: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (2019)

[62] Huang, B.: [03/24] Boston Housing Dataset. https://kaggle.com/competitions/2403-boston-housing-dataset. Kaggle (2020)

[63] Kandhro, I.A., Jumani, S.Z., Ali, F., Shaikh, Z.U., Arain, M.A., Shaikh, A.A.: Performance analysis of hyperparameters on a sentiment analysis model. Engineering, Technology & Applied Science Research **10**(4), 6016–6020 (2020)

[64] Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. arXiv preprint arXiv:1301.3781 (2013)

[65] Orús, R.: A practical introduction to tensor networks: Matrix product states and projected entangled pair states. Annals of physics **349**, 117–158 (2014)

[66] Li, Z., Huang, C., Wang, X., Hu, H., Wyeth, C., Bu, D., Yu, Q., Gao, W., Liu, X., Li, M.: Lossless data compression by large models. Nature Machine Intelligence, 1–6 (2025)

[67] Berezutskii, A., Liu, M., Acharya, A., Ellerbrock, R., Gray, J., Haghshenas, R., He, Z., Khan, A., Kuzmin, V., Lyakh, D., et al.: Tensor networks for quantum computing. Nature Reviews Physics, 1–13 (2025)

[68] Sitzmann, V., Martel, J., Bergman, A., Lindell, D., Wetzstein, G.: Implicit neural representations with periodic activation functions. Advances in neural information processing systems **33**, 7462–7473 (2020)

[69] He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 16000–16009 (2022)

[70] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778 (2016)

[71] Zimmer, L., Lindauer, M., Hutter, F.: Auto-pytorch: Multi-fidelity metalearning for efficient and robust autodl. IEEE transactions on pattern analysis and machine intelligence **43**(9), 3079–3090 (2021)

[72] Isensee, F., Jaeger, P.F., Kohl, S.A., Petersen, J., Maier-Hein, K.H.: nnu-net: a self-configuring method for deep learning-based biomedical image segmentation.

Nature methods **18**(2), 203–211 (2021)

[73] Isensee, F., Wald, T., Ulrich, C., Baumgartner, M., Roy, S., Maier-Hein, K., Jaeger, P.F.: nnu-net revisited: A call for rigorous validation in 3d medical image segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, pp. 488–498 (2024). Springer

[74] Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical Image Computing and Computer-assisted Intervention, pp. 234–241 (2015). Springer

[75] Isensee, F., Jaeger, P.F., Kohl, S.A., Petersen, J., Maier-Hein, K.H.: nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. Nature methods **18**(2), 203–211 (2021)

[76] Roy, S., Koehler, G., Ulrich, C., Baumgartner, M., Petersen, J., Isensee, F., Jaeger, P.F., Maier-Hein, K.H.: Mednext: transformer-driven scaling of convnets for medical image segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, pp. 405–415 (2023). Springer

[77] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017)

[78] Tang, Y., Yang, D., Li, W., Roth, H.R., Landman, B., Xu, D., Nath, V., Hatamizadeh, A.: Self-supervised pre-training of swin transformers for 3d medical image analysis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 20730–20740 (2022)

[79] Gu, A., Dao, T.: Mamba: Linear-time sequence modeling with selective state spaces. arXiv preprint arXiv:2312.00752 (2023)

[80] Ma, J., Li, F., Wang, B.: U-mamba: Enhancing long-range dependency for biomedical image segmentation. arXiv preprint arXiv:2401.04722 (2024)

[81] Landman, B., Xu, Z., Igelsias, J., Styner, M., Langerak, T., Klein, A.: Miccai multi-atlas labeling beyond the cranial vault–workshop and challenge. In: Proc. MICCAI Multi-Atlas Labeling Beyond Cranial Vault—Workshop Challenge, vol. 5, p. 12 (2015)

[82] LaBella, D., Adewole, M., Alonso-Basanta, M., Altes, T., Anwar, S.M., Baid, U., Bergquist, T., Bhalerao, R., Chen, S., Chung, V., Conte, G.-M., Dako, F., Eddy, J., Ezhov, I., Godfrey, D., Hilal, F., Familiar, A., Farahani, K., Iglesias, J.E., Jiang, Z., Johanson, E., Kazerooni, A.F., Kent, C., Kirkpatrick, J., Kofler, F., Leemput, K.V., Li, H.B., Liu, X., Mahtabfar, A., McBurney-Lin, S., McLean, R., Meier, Z., Moawad, A.W., Mongan, J., Nedelec, P., Pajot, M., Piraud, M., Rashid, A., Reitman, Z., Shinohara, R.T., Velichko, Y., Wang, C., Warman, P.,

Wiggins, W., Aboian, M., Albrecht, J., Anazodo, U., Bakas, S., Flanders, A., Janas, A., Khanna, G., Linguraru, M.G., Menze, B., Nada, A., Rauschecker, A.M., Rudie, J., Tahon, N.H., Villanueva-Meyer, J., Wiestler, B., Calabrese, E.: The ASNR-MICCAI Brain Tumor Segmentation (BraTS) Challenge 2023: Intracranial Meningioma (2023)

[83] Jiang, H., Wang, J., Cong, W., Huang, Y., Ramezani, M., Sarma, A., Dokholyan, N.V., Mahdavi, M., Kandemir, M.T.: Predicting protein–ligand docking structure with graph neural network. Journal of chemical information and modeling **62**(12), 2923–2932 (2022)

[84] Jiang, D., Hsieh, C.-Y., Wu, Z., Kang, Y., Wang, J., Wang, E., Liao, B., Shen, C., Xu, L., Wu, J., *et al.*: Interactiongraphnet: a novel and efficient deep graph representation learning framework for accurate protein–ligand interaction predictions. Journal of medicinal chemistry **64**(24), 18209–18232 (2021)

[85] Shi, W., Yang, H., Xie, L., Yin, X.-X., Zhang, Y.: A review of machine learning-based methods for predicting drug–target interactions. Health Information Science and Systems **12**(1), 30 (2024)

[86] Gu, S., Bao, L., Yang, Y., Zhao, Y., Tong, H.H.Y., Liu, L., Liu, H., Hou, T., Kang, Y.: Amgc is a multiple-task graph neutral network for epigenetic target profiling. Cell Reports Physical Science **5**(3) (2024)

[87] Bongini, P., Bianchini, M., Scarselli, F.: Molecular generative graph neural networks for drug discovery. Neurocomputing **450**, 242–252 (2021)

[88] Wang, X., Ma, Y., Wang, Y., Jin, W., Wang, X., Tang, J., Jia, C., Yu, J.: Traffic flow prediction via spatial temporal graph neural network. In: Proceedings of the Web Conference 2020, pp. 1082–1092 (2020)

[89] Sharma, A., Sharma, A., Nikashina, P., Gavrilenko, V., Tselykh, A., Bozhenyuk, A., Masud, M., Meshref, H.: A graph neural network (gnn)-based approach for real-time estimation of traffic speed in sustainable smart cities. Sustainability **15**(15), 11893 (2023)

[90] Mirhoseini, A., Goldie, A., Yazgan, M., Jiang, J.W., Songhori, E., Wang, S., Lee, Y.-J., Johnson, E., Pathak, O., Nova, A., *et al.*: A graph placement methodology for fast chip design. Nature **594**(7862), 207–212 (2021)

[91] Yang, S., Yang, Z., Li, D., Zhang, Y., Zhang, Z., Song, G., Hao, J.: Versatile multi-stage graph neural network for circuit representation. Advances in Neural Information Processing Systems **35**, 20313–20324 (2022)

[92] Zhao, T., Zhang, X., Wang, S.: Graphsmote: Imbalanced node classification on graphs with graph neural networks. In: Proceedings of the 14th ACM International Conference on Web Search and Data Mining, pp. 833–841 (2021)

[93] Yan, L., Zhang, S., Li, B., Zhou, M., Huang, Z.: Unreal: Unlabeled nodes retrieval and labeling for heavily-imbalanced node classification. arXiv preprint arXiv:2303.10371 (2023)

[94] Korotcov, A., Tkachenko, V., Russo, D.P., Ekins, S.: Comparison of deep learning with multiple machine learning methods and metrics using diverse drug discovery data sets. Molecular pharmaceutics **14**(12), 4462–4475 (2017)

[95] Ekins, S., Siqueira-Neto, J., McCall, L.-I., Sarker, M., Yadav, M., Ponder, E.L., Kallel, E.A., Kellar, D., Chen, S., Arkin, M., *et al.*: Machine learning models and pathway genome data base for trypanosoma cruzi drug discovery. PLoS neglected tropical diseases **9**(6), 0003878 (2015)

[96] Weininger, D.: Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. Journal of chemical information and computer sciences **28**(1), 31–36 (1988)

[97] Duvenaud, D.K., Maclaurin, D., Iparraguirre, J., Bombarell, R., Hirzel, T., Aspuru-Guzik, A., Adams, R.P.: Convolutional networks on graphs for learning molecular fingerprints. Advances in neural information processing systems **28** (2015)

[98] Fliri, A.F., Loging, W.T., Thadeio, P.F., Volkmann, R.A.: Biological spectra analysis: linking biological activity profiles to molecular structure. Proceedings of the National Academy of Sciences **102**(2), 261–266 (2005)

[99] Chan, J.S., Chowdhury, N., Jaffe, O., Aung, J., Sherburn, D., Mays, E., Starace, G., Liu, K., Maksin, L., Patwardhan, T., et al.: Mle-bench: Evaluating machine learning agents on machine learning engineering. arXiv preprint arXiv:2410.07095 (2024)

[100] Hill, T.P.: Knowing when to stop: How to gamble if you must-the mathematics of optimal stopping. American Scientist **97**(2), 126–133 (2009)

[101] Ferguson, T.S.: Who solved the secretary problem? Statistical science **4**(3), 282–289 (1989)