# Language Modeling by Language Models

Junyan Cheng, Peter Clark, **Kyle Richardson**

May 2025

**Allen Institute for Artificial Intelligence** (AI2)

# The big picture

- **Fully autonomous discovery**, simulate all aspects of the conventional research process *(e.g., ideation, experiment execution, paper writing)*.

# The big picture

- **Autonomous discovery for ML**: Discovering novel machine learning components, make our ML systems more efficient, transparent and safer..

# The big picture

▶ **Autonomous discovery for ML**: Discovering novel machine learning components, make our ML systems more efficient, transparent and safer..

Lion optimizer (Chen et al., 2023)

# The big picture

▶ **Autonomous discovery for ML**: Discovering novel machine learning components, make our ML systems more efficient, transparent and safer..

Lion optimizer (Chen et al., 2023)



**Symbolic Discovery of Optimization Algorithms**

Narrow (no LLMs), clear goals.

# The big picture

▶ **Autonomous discovery for ML**: Discovering novel machine learning components, make our ML systems more efficient, transparent and safer..

Lion optimizer (Chen et al., 2023)



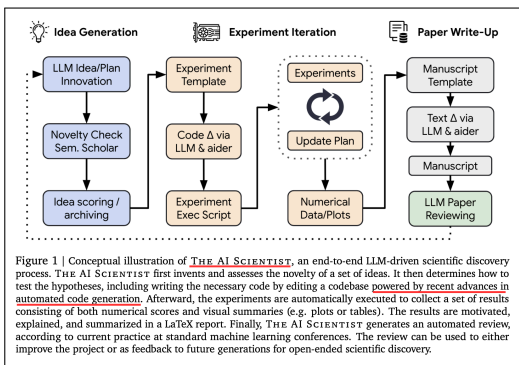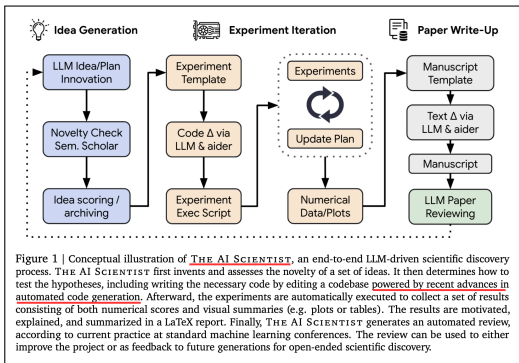**Symbolic Discovery of Optimization Algorithms**

Xiangning Chen[1 2 1 *]    Chen Liang[1 1]    Da Huang[1]    Esteban Real[1]

Kaiyuan Wang[1]    Hieu Pham[1]    Xuanyi Dong[1]    Thang Luong[1]

Cho-Jui Hsieh[2]    Yifeng Lu[1]    Quoc V. Le[1]

[1]Equal & Core Contribution

[1]Google          [2]UCLA

**Abstract**

We present a method to formulate algorithm discovery as program search, and apply it to discover optimization algorithms for deep neural network training. We leverage efficient search techniques to explore an infinite and sparse program space. To bridge the large generalization gap between proxy and target tasks, we also introduce program selection and simplification strategies. Our method discovers a simple and effective optimization algorithm, **Lion** (*EvoLved Sign Momentum*). It is more memory-efficient than Adam as it only keeps track of the momentum. Different from adaptive optimizers, its update has the same magnitude for each parameter calculated through the sign operation. We compare Lion with widely used optimizers, such as Adam and Adafactor, for training a variety of models on different tasks. On image classification, Lion boosts the accuracy of ViT by up to 2% on ImageNet and saves up to 5x the pre-training compute on JFT. On vision-language contrastive learning, we achieve 88.3% *zero-shot* and 91.1% *fine-tuning* accuracy on ImageNet, surpassing the previous best results by 2% and 0.1%, respectively. On diffusion models, Lion outperforms Adam by achieving a better FID score and reducing the training compute by up to 2.3x. For autoregressive, masked language modeling, and fine-tuning, Lion exhibits a similar or better performance compared to Adam. Our analysis of Lion reveals that its performance gain grows with the training batch size. It also requires a smaller learning rate than Adam due to the larger norm of the update produced by the sign function. Additionally, we examine the limitations of Lion and identify scenarios where its improvements are small or not statistically significant.

Narrow (no LLMs), clear goals.

AI Scientist (Lu et al., 2024)



Figure 1 | Conceptual illustration of THE AI SCIENTIST, an end-to-end LLM-driven scientific discovery process. THE AI SCIENTIST first invents and assesses the novelty of a set of ideas. It then determines how to test the hypotheses, including writing the necessary code by editing a codebase powered by recent advances in automated code generation. Afterward, the experiments are automatically executed to collect a set of results consisting of both numerical scores and visual summaries (e.g. plots or tables). The results are motivated, explained, and summarized in a LaTeX report. Finally, THE AI SCIENTIST generates an automated review, according to current practice at standard machine learning conferences. The review can be used to either improve the project or as feedback to future generations for open-ended scientific discovery.

# The big picture

▶ **Autonomous discovery for ML**: Discovering novel machine learning components, make our ML systems more efficient, transparent and safer..

Lion optimizer (Chen et al., 2023)



Narrow (no LLMs), clear goals.

AI Scientist (Lu et al., 2024)



Figure 1 | Conceptual illustration of THE AI SCIENTIST, an end-to-end LLM-driven scientific discovery process. THE AI SCIENTIST first invents and assesses the novelty of a set of ideas. It then determines how to test the hypotheses, including writing the necessary code by editing a codebase powered by recent advances in automated code generation. Afterward, the experiments are automatically executed to collect a set of results consisting of both numerical scores and visual summaries (e.g. plots or tables). The results are motivated, explained, and summarized in a LaTeX report. Finally, THE AI SCIENTIST generates an automated review, according to current practice at standard machine learning conferences. The review can be used to either improve the project or as feedback to future generations for open-ended scientific discovery.

Broad (LLM-driven), unclear goals.

Lion optimizer (Chen et al., 2023)

AI Scientist (Lu et al., 2024)

**Problem**: Largely system demonstrations

Narrow (no LLMs), clear goals.

Broad (LLM-driven), unclear goals.

2

# 1. **Tasks:** What are the target discovery tasks?



Figure 1 | Conceptual illustration of <u>THE AI SCIENTIST</u>, an end-to-end LLM-driven scientific discovery process. THE AI SCIENTIST first invents and assesses the novelty of a set of ideas. It then determines how to test the hypotheses, including writing the necessary code by editing a codebase <u>powered by recent advances in automated code generation</u>. Afterward, the experiments are automatically executed to collect a set of results consisting of both numerical scores and visual summaries (e.g. plots or tables). The results are motivated, explained, and summarized in a LaTeX report. Finally, THE AI SCIENTIST generates an automated review, according to current practice at standard machine learning conferences. The review can be used to either improve the project or as feedback to future generations for open-ended scientific discovery.

▶ What tasks and discovery problems should we be working on to make progress? **Community has not yet come up with clear tasks or metrics.**

# 1. **Tasks:** What are the target discovery tasks?

Lion optimizer (Chen et al., 2023)



Narrow (no LLMs), <u>clear goals</u>.

AI Scientist (Lu et al., 2024)
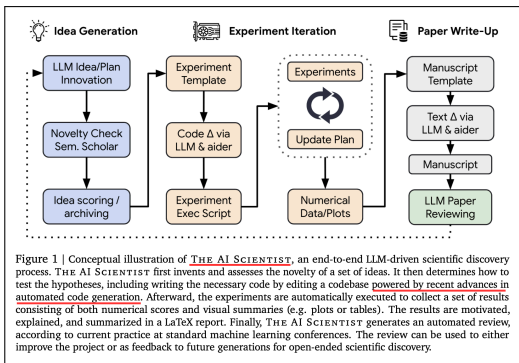


Figure 1 | Conceptual illustration of THE AI SCIENTIST, an end-to-end LLM-driven scientific discovery process. THE AI SCIENTIST first invents and assesses the novelty of a set of ideas. It then determines how to test the hypotheses, including writing the necessary code by editing a codebase powered by recent advances in automated code generation. Afterward, the experiments are automatically executed to collect a set of results consisting of both numerical scores and visual summaries (e.g. plots or tables). The results are motivated, explained, and summarized in a LaTeX report. Finally, THE AI SCIENTIST generates an automated review, according to current practice at standard machine learning conferences. The review can be used to either improve the project or as feedback to future generations for open-ended scientific discovery.

<u>Broad</u> (LLM-driven), unclear goals.

# 1. **Tasks:** What are the target discovery tasks?

Lion optimizer (Chen et al., 2023)



Narrow (no LLMs), <u>clear goals</u>.

AI Scientist (Lu et al., 2024)



Figure 1 | Conceptual illustration of THE AI SCIENTIST, an end-to-end LLM-driven scientific discovery process. THE AI SCIENTIST first invents and assesses the novelty of a set of ideas. It then determines how to test the hypotheses, including writing the necessary code by editing a codebase powered by recent advances in automated code generation. Afterward, the experiments are automatically executed to collect a set of results consisting of both numerical scores and visual summaries (e.g. plots or tables). The results are motivated, explained, and summarized in a LaTeX report. Finally, THE AI SCIENTIST generates an automated review, according to current practice at standard machine learning conferences. The review can be used to either improve the project or as feedback to future generations for open-ended scientific discovery.

<u>Broad</u> (LLM-driven), unclear goals.

**Our proposal**: language model architecture discovery, finding better (e.g., more efficient, performant, transparent,…), LM layer designs.

3

## 2. **System design**: How should we build such systems?



Figure 1 | Conceptual illustration of THE AI SCIENTIST, an end-to-end LLM-driven scientific discovery process. THE AI SCIENTIST first invents and assesses the novelty of a set of ideas. It then determines how to test the hypotheses, including writing the necessary code by editing a codebase powered by recent advances in automated code generation. Afterward, the experiments are automatically executed to collect a set of results consisting of both numerical scores and visual summaries (e.g. plots or tables). The results are motivated, explained, and summarized in a LaTeX report. Finally, THE AI SCIENTIST generates an automated review, according to current practice at standard machine learning conferences. The review can be used to either improve the project or as feedback to future generations for open-ended scientific discovery.

▶ Do these discovery workflows make sense? What are their limits? Should be efficient, cost-effective, transparent.

# 2. **System design**: How should we build such systems?



Figure 1 | Conceptual illustration of THE AI SCIENTIST, an end-to-end LLM-driven scientific discovery process. THE AI SCIENTIST first invents and assesses the novelty of a set of ideas. It then determines how to test the hypotheses, including writing the necessary code by editing a codebase powered by recent advances in automated code generation. Afterward, the experiments are automatically executed to collect a set of results consisting of both numerical scores and visual summaries (e.g. plots or tables). The results are motivated, explained, and summarized in a LaTeX report. Finally, THE AI SCIENTIST generates an automated review, according to current practice at standard machine learning conferences. The review can be used to either improve the project or as feedback to future generations for open-ended scientific discovery.

## 2. **System design**: How should we build such systems?



Figure 1 | Conceptual illustration of THE AI SCIENTIST, an end-to-end LLM-driven scientific discovery process. THE AI SCIENTIST first invents and assesses the novelty of a set of ideas. It then determines how to test the hypotheses, including writing the necessary code by editing a codebase powered by recent advances in automated code generation. Afterward, the experiments are automatically executed to collect a set of results consisting of both numerical scores and visual summaries (e.g. plots or tables). The results are motivated, explained, and summarized in a LaTeX report. Finally, THE AI SCIENTIST generates an automated review, according to current practice at standard machine learning conferences. The review can be used to either improve the project or as feedback to future generations for open-ended scientific discovery.

**Proposed** a new algorithmic framework for discovery, allows us to address technical issues, devise generalized algorithms.

Figure 1 | Conceptual illustration of THE AI SCIENTIST, an end-to-end LLM-driven scientific discovery process. THE AI SCIENTIST first invents and assesses the novelty of a set of ideas. It then determines how to test the hypotheses, including writing the necessary code by editing a codebase powered by recent advances in automated code generation. Afterward, the experiments are automatically executed to collect a set of results consisting of both numerical scores and visual summaries (e.g. plots or tables). The results are motivated, explained, and summarized in a LaTeX report. Finally, THE AI SCIENTIST generates an automated review, according to current practice at standard machine learning conferences. The review can be used to either improve the project or as feedback to future generations for open-ended scientific discovery.

**Problem**: Language model architecture discovery

**Proposed** a new algorithmic framework for discovery, allows us to address technical issues, devise generalized algorithms.

# Language model architecture design discovery: what?

- ▶ Finding improved layer designs for autoregressive language models.

# Language model architecture design discovery: what?

- ▶ Finding improved layer designs for autoregressive language models.

# Language model architecture design discovery: what?

▶ Finding improved layer designs for autoregressive language models.



```
def GPT(X,**Z):
    X1,Z = RMSNorm(X,**Z)
    X2,Z = MHA(X1,**Z)
    X = X+X2
    X3,Z = RMSNorm(X,**Z)
    X4,Z = GatedMLP(X3,**Z)
    X = X+X4
    return X,Z
```

# Language model architecture design discovery: what?

▶ Finding improved layer designs for autoregressive language models.



Can we improve on these block designs?

At its core, a **code discovery** problem, similar goals to AutoML and Neural architecture search (**NAS**), model full research pipeline.

# Why is this an interesting problem?



TechCrunch

Latest  Startups  Venture  Apple  Security  AI  Apps  Startup Battlefield  |  Events  Podcasts  Newsletters    Sign In

AI

## TTT models might be the next frontier in generative AI

Kyle Wiggers · 1:47 PM PDT · July 17, 2024

IMAGE CREDITS: NATEE127 / GETTY IMAGES

After years of dominance by the form of AI known as the transformer, the hunt is on for new architectures.

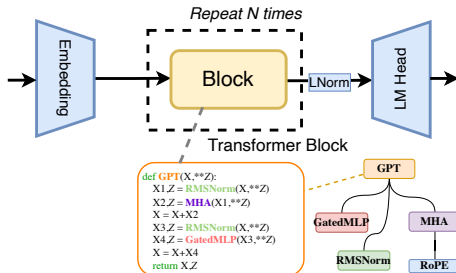Transformers underpin OpenAI's video-generating model Sora, and they're at the heart of text-generating models like Anthropic's Claude, Google's Gemini and GPT-4o. But they're beginning to run up against technical roadblocks — in particular, computation-related roadblocks.

# Why is this an interesting problem?



Retentive Network: A Successor to Transformer
for Large Language Models

Language Models are Unsupervised Multitask Learners

# Why is this an interesting problem?

▶ Finding improved layer designs for autoregressive language models.



**Ill-formed search space**: huge unbounded design space.

**Complex sampling process**: literature understanding, coding skills.

**Expensive verification**: pre-training/evaluation, resource bound.

# Language model architecture design discovery: how?

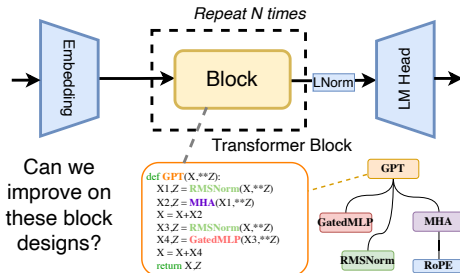▶ Finding improved layer designs for auto-regressive language models.



**Continuous learning loop**: Generate new model ideas, implement them and verify through generative pre-training.

# Language model architecture design discovery: how?

▶ Finding improved layer designs for auto-regressive language models.



**Continuous learning loop**: Generate new model ideas, implement them and verify through generative pre-training.

▶ **Objective**: Find designs that improve on end-task performance.

# Language model architecture design discovery: how?

▶ Finding improved layer designs for auto-regressive language models.
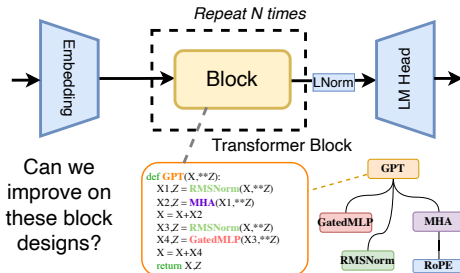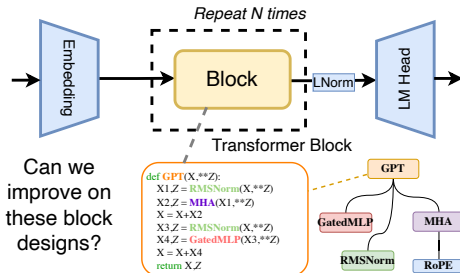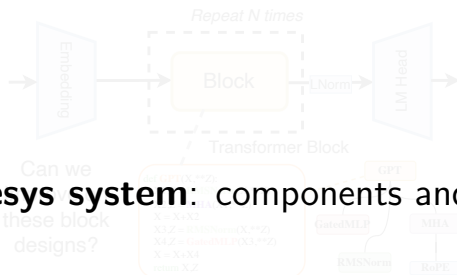


**Continuous learning loop**: Generate new model ideas, implement them and verify through generative pre-training.

▶ **Objective**: Find designs that improve on end-task performance.
▶ Start small, innovate then scale, **Ladder-of-scales** (LoS) approach.

# The **Genesys system**: components and principles

# The Genesys system: **core utilities**



library

External Sources

Semantic **Scholar**

ar𝕏iv

Papers With Code

Paper Vector DB

❀Pinecone  **mathpix**

Reference Library

Web Search Agent

⊗ perplexity

verifier

Automated Trainer

🦎SmolLM

🤗 **Hugging Face**

Weights & Biases

Realtime Leaderboard

Automated Evaluator

◉ **LM-Eval** BLiMP

⌁ GLUE

tree

# The Genesys system: **agents**

# The Genesys system: **agents**

# The Genesys system: **distributed evolution**



## library

**External Sources**

Semantic **Scholar**

arXiv

Papers With Code

**Paper Vector DB**

Pinecone    mathpix

**Reference Library**

**Web Search Agent**

perplexity

## verifier

**Automated Trainer**

SmolLM

🤗 **Hugging Face**

Weights & Biases

**Realtime Leaderboard**

**Automated Evaluator**

LM-Eval   BLiMP   ...   GLUE

design
**Designer**
*Parent Designs*
**Selector**

**Experimenter**
*Design & Scale*
**Verifier**
*verify*

*Responses*
*Queries*

*Design & Scale*
*Results*

*New Design*
*Report*

## evolution

```
async def genesys_evolve(EvoTree, KE, VE):
    """Run distributed discovery loop"""
    while budget:
        async design(EvoTree,KE)
        async verify(EvoTree,VE)
```
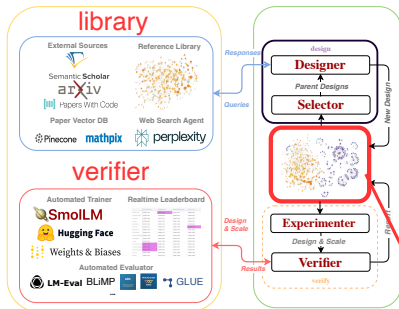
# The Genesys system



**Experiments at a glance**: 1,162 discovered designs (1,062 fully verified), 86K dialogues, 2.76M lines of code, 1B processed tokens.
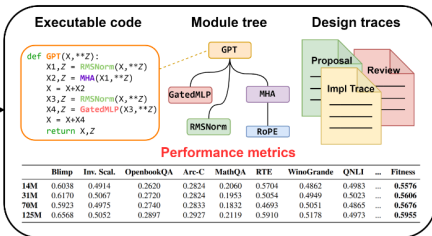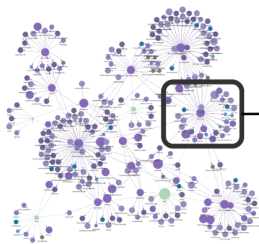
Peeking inside

**Experiments at a glance**: 1,162 discovered designs (1,062 fully verified), 86K dialogues, 2.76M lines of code, 1B processed tokens.
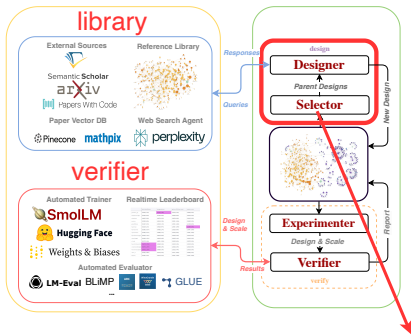
# Design tree: fully factorizable design space



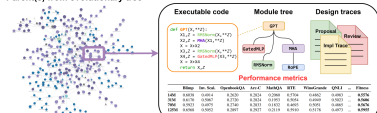Code is fully factorizable, representable as a unit tree

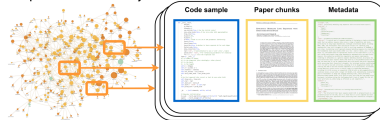**Fitness score:** end task performance

library

External Sources / Reference Library

Semantic Scholar
arXiv
Papers With Code

Paper Vector DB / Web Search Agent
Pinecone / mathpix / perplexity

verifier

Automated Trainer / Realtime Leaderboard
SmolLM
Hugging Face
Weights & Biases

Automated Evaluator
LM-Eval BLiMP / GLUE

Designer
Selector
Experimenter
Verifier

Executable code

```
def GPT(X,**Z):
    X1,Z = RMSNorm(X,**Z)
    X2,Z = MHA(X1,**Z)
    X = X+X2
    X3,Z = RMSNorm(X,**Z)
    X4,Z = GatedMLP(X3,**Z)
    X = X+X4
    return X,Z
```

Module tree

GPT
GatedMLP — MHA
RMSNorm — RoPE

Design traces

Proposal / Review
Impl Trace

## Performance metrics

|       | Blimp  | Inv. Scal. | OpenbookQA | Arc-C  | MathQA | RTE    | WinoGrande | QNLI   | ... | Fitness |
|-------|--------|------------|------------|--------|--------|--------|------------|--------|-----|---------|
| 14M   | 0.6038 | 0.4914     | 0.2620     | 0.2824 | 0.2060 | 0.5704 | 0.4862     | 0.4983 | ... | 0.5576  |
| 31M   | 0.6170 | 0.5067     | 0.2720     | 0.2824 | 0.1953 | 0.5054 | 0.4949     | 0.5023 | ... | 0.5606  |
| 70M   | 0.5923 | 0.4975     | 0.2740     | 0.2833 | 0.1832 | 0.4693 | 0.5051     | 0.4865 | ... | 0.5676  |
| 125M  | 0.6568 | 0.5052     | 0.2897     | 0.2927 | 0.2119 | 0.5910 | 0.5178     | 0.4973 | ... | 0.5955  |

# **Designers**: Proposer-reviewer architecture

# **Designers**: Planner, coder, observer



Unit-based design and implementation

# **Verifiers**: budget sensitive scaling



Ladder of scales (LoS)
approach

Have we made any discoveries yet?

# A sketch of the results: **end task performance**

| | **Blimp** | **Wnli** | **RTE** | **WG** | **CoLA** | **SST2** | **WSC** | **IS** | **Mrpc** | **avg.** |
|---|---|---|---|---|---|---|---|---|---|---|
| *Random* | 69.75 | 43.66 | 52.71 | 48.78 | 50.00 | 49.08 | 49.82 | 50.03 | 31.62 | 49.49 |
| GPT | <u>92.70</u> | 60.56 | <u>62.80</u> | 52.17 | *53.24* | 54.13 | 56.76 | 55.31 | 68.38 | <u>61.78</u> |
| Mamba2 | 83.22 | 63.38 | **63.88** | 51.22 | <u>55.94</u> | <u>56.58</u> | <u>57.12</u> | 53.85 | 67.89 | 61.45 |
| RWKV7 | 88.76 | 61.97 | 60.21 | *49.80* | 54.25 | 55.32 | *54.57* | **57.00** | 68.38 | 61.14 |
| RetNet | 85.16 | 61.97 | 61.35 | 50.51 | **56.29** | 55.43 | 56.03 | 54.95 | *56.37* | 59.78 |
| TTT | 86.13 | 63.38 | *55.23* | 50.75 | 55.55 | 56.35 | 54.93 | 55.31 | 59.80 | 59.71 |
| VQH | **94.37** | 59.15 | 59.91 | 50.28 | 54.25 | 53.56 | *53.83* | 49.45 | 56.62 | 59.05 |
| HMamba | 83.74 | <u>64.79</u> | 61.35 | **53.59** | 54.69 | **57.04** | 56.40 | 54.58 | 59.31 | 60.61 |
| Geogate | 90.95 | <u>59.15</u> | 61.35 | <u>52.72</u> | 54.25 | 55.32 | **58.96** | 54.95 | <u>68.63</u> | **61.81** |
| Hippovq | 87.96 | *50.70* | 59.91 | <u>50.28</u> | 54.25 | 55.73 | *53.83* | <u>55.68</u> | **69.88** | 59.80 |
| SRN | *80.83* | **65.52** | 59.55 | 50.75 | 54.45 | *52.98* | 56.03 | <u>54.95</u> | 61.03 | *59.57* |

Table 3: Performance of human designs and discovered models on various Benchmarks (350M Parameters, 50B Tokens). Metrics indicate accuracy percentages. Bold and underlined denotes the top and second best, italics denoting worst.

**new designs**

|  | Blimp | Wnli | RTE | WG | CoLA | SST2 | WSC | IS | Mrpc | avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| *Random* | 69.75 | 43.66 | 52.71 | 48.78 | 50.00 | 49.08 | 49.82 | 50.03 | 31.62 | 49.49 |
| GPT | 92.70 | 60.56 | 62.80 | 52.17 | *53.24* | 54.13 | 56.76 | 55.31 | 68.38 | 61.78 |
| Mamba2 | 83.22 | 63.38 | **63.88** | 51.22 | 55.94 | 56.58 | 57.12 | 53.85 | 67.89 | 61.45 |
| RWKV7 | 88.76 | 61.97 | 60.21 | *49.80* | 54.25 | 55.32 | 54.57 | **57.00** | 68.38 | 61.14 |
| RetNet | 85.16 | 61.97 | 61.35 | 50.51 | **56.29** | 55.43 | 56.03 | 54.95 | *56.37* | 59.78 |
| TTT | 86.13 | 63.38 | *55.23* | 50.75 | 55.55 | 56.35 | 54.93 | 55.31 | 59.80 | 59.71 |
| VO | | | | | | | | | | |
| Geogate | 90.95 | 59.15 | 61.35 | 52.72 | 54.25 | 55.32 | **58.96** | 54.95 | 68.63 | **61.81** |
| Hippovq | 87.96 | *50.70* | 59.91 | 50.28 | 54.25 | 55.73 | *53.83* | 55.68 | **69.88** | 59.80 |
| SRN | *80.83* | **65.52** | 59.55 | 50.75 | 54.45 | *52.98* | 56.03 | 54.95 | 61.03 | *59.57* |

Table 3: Performance of human designs and discovered models on various Benchmarks (350M Parameters, 50B Tokens). Metrics indicate accuracy percentages. Bold and underlined denotes the top and second best, italics denoting worst.

new designs

**Result**: Yields designs competitive with human ones

8

# A sketch of the results: **system and design analysis**



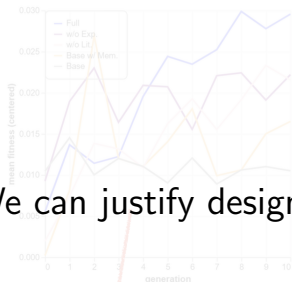Table 3. Agent benchmark results. Bold and underlined denotes the top and second best. "Library" stands for our reference library with 180 designs providing core block code.

| | Valid | Attempts | Costs | LFC |
|---|---|---|---|---|
| Full | 92% | 2.6 (±1.1) | 15.0 (±18.5) | 181 (±44) |
| No FF | 73% | 3.0 (±1.7) | 7.9 (±7.1) | 75 (±29) |
| No Pl. | 91% | 2.6 (±1.1) | 16.0 (±20.8) | 218 (±69) |
| No Ob. | 89% | 2.6 (±1.1) | 12.1 (±20.1) | 211 (±67) |
| No SC | 30% | 2.4 (±1.0) | 2.9 (±4.7) | 167 (±33) |
| Simple | 6% | 1.1 (±0.2) | 0.3 (±0.3) | 49 (±15) |
| *Library* | - | - | - | 220 (±136) |

**system stability**

**successful code generation rates**

We can justify design, empirically and formally.

system stability

successful code
generation rates

Please come to the poster to learn more

Thank you.

# References I

Chen, X., Liang, C., Huang, D., Real, E., Wang, K., Pham, H., Dong, X., Luong, T., Hsieh, C.-J., Lu, Y., et al. (2023). Symbolic discovery of optimization algorithms. *Advances in neural information processing systems*, 36:49205–49233.

Lu, C., Lu, C., Lange, R. T., Foerster, J., Clune, J., and Ha, D. (2024). The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*.