# AI as a Scientific Pilot ?

Karteek Alahari          Christophe Biernacki

Prepare yourself to become an AI-assisted researcher

(or change your job !)

# Have we been here before ?

# Outline

- **Introduction**

- A focus on AI Scientist-v2 – main principles

- Investigating the next stages (foresight)

- Conclusion

# Impressive progress – AI Tools for Academia

- Finding, analyzing, summarizing academic articles: PaperPilot, AIModels.fyi

- Assistance for writing: Paperpal

- Measuring compliance of open science requirements: Dataseer

- Evaluating AI's ability to replicate AI research: PaperBench

PaperPilot https://www.paperpilot.xyz/
AIModels.fyi https://www.aimodels.fyi/
Paperpal https://paperpal.com
Dataseer https://dataseer.ai
PaperBench https://arxiv.org/abs/2504.01848

# Even more impressive – Stages beyond

- Finding potential trends: NotebookLM

- Novel hypotheses / research plans: AI co-scientist, Virtual Scientists (VirSci)

- Generating research papers
  - with some human intervention (e.g., manuscript preparation): Zochi
  - more autonomously: AI Scientist-v2
  - **One such paper accepted at an ICLR 2025 workshop***

NotebookLM https://notebooklm.google
AI co-scientist https://arxiv.org/abs/2502.18864
Virtual Scientists https://arxiv.org/abs/2410.09403
AI Scientist-v2 https://arxiv.org/abs/2504.08066
Zochi https://github.com/IntologyAI/Zochi
* https://github.com/SakanaAI/AI-Scientist-ICLR2025-Workshop-Experiment

# Even more impressive – Stages beyond

- Finding potential trends: NotebookLM

- Novel hypotheses / research plans: AI co-scientist, Virtual Scientists (VirSci)

## Can AI systems take the role of principal investigator (PI) in research?

- Generating research papers autonomously: AI Scientist-v2

- **One such paper accepted at an ICLR 2025 workshop***

# Outline

- Introduction

- **A focus on AI Scientist-v2 – main principles**

- Investigating the next stages (foresight)

- Conclusion

# 3 AI Generated Papers

| Title | Workshop result |
| --- | --- |
| **Compositional Regularization: Unexpected Obstacles in Enhancing Neural Network Generalization** | Accepted (scores 6, 7, 6) |
| Real-World Challenges in Pest Detection Using Deep Learning: An Investigation into Failures and Solutions | Rejected (scores: 3, 7, 4) |
| Unveiling the Impact of Label Noise on Model Calibration in Deep Learning | Rejected (scores: 3, 3, 3) |

- Total time for generation: several to 15 hours (fixed runtime limit)

- The "accepted" manuscript
  - ~ top 45% of submissions

  - first fully AI-generated manuscript to successfully pass a peer-review process

# 3 AI Generated Papers

- Entirely generated end-to-end by AI - no modifications from humans

- AI Scientist-v2 came up with the scientific hypothesis

- Proposed experiments to test the hypothesis

- Wrote and evaluated code to conduct the experiments

- Ran the experiments, analyzed the data, produced figures

- Wrote the entire document

- **Only human input**: broad topic of research (to be relevant to the workshop)

# The "Accepted" paper



- respects the paper format of the workshop (length, layout, references…)

- includes: formulas, figures, references, experiments, supplementary material…

# AI Scientist-v2 Workflow



- **Key technical points**: agentic tree search, VLM (Vision Language Models), feedback, and parallel experiment execution
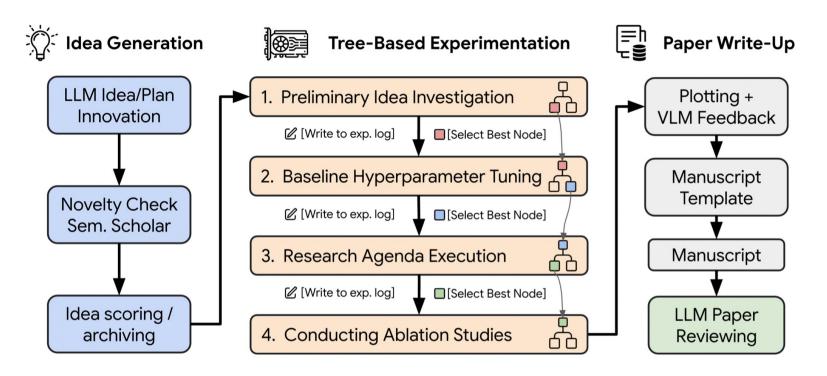
# Interaction with Humans: Prompts

## Idea Generation Prompt
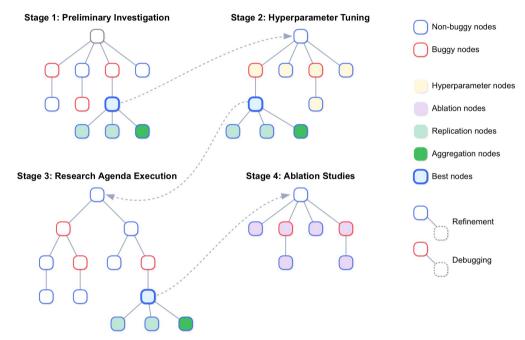
```
# System prompt
You are an experienced AI researcher who aims to propose high-impact
research ideas resembling exciting grant proposals. Feel free to propose
any novel ideas or experiments; make sure they are novel. Be very creative
and think out of the box. Each proposal should stem from a simple and
elegant question, observation, or hypothesis about the topic. For example,
they could involve very interesting and simple interventions or
investigations that explore new possibilities or challenge existing
assumptions. Clearly clarify how the proposal distinguishes from
the existing literature.

Ensure that the proposal can be done starting from the provided
codebase, and does not require resources beyond what an academic
lab could afford. These proposals should lead to papers that are
publishable at top ML conferences.
```

```
You have access to the following tools:

{tool_descriptions}

Respond in the following format:

ACTION:
<The action to take, exactly one of {tool_names_str}>

ARGUMENTS:
<If ACTION is "SearchSemanticScholar", provide the search query
as {{"query": "your search query"}}. If ACTION is "FinalizeIdea",
provide the idea details as {{"idea": {{ ... }}}} with the IDEA JSON
specified below.>

If you choose to finalize your idea, provide the IDEA JSON in the arguments:

IDEA JSON:
```json
{{
    "Name": "...",
    "Title": "...",
    "Short Hypothesis": "...",
    "Related Work": "...",
    "Abstract": "...",
    "Experiments": "...",
    "Risk Factors and Limitations": "..."
}}
```

Ensure the JSON is properly formatted for automatic parsing.

Note: You should perform at least one literature search before finalizing
your idea to ensure it is well-informed by existing research.

# Initial idea generation prompt
{workshop_description}

Here are the proposals that you have already generated:

{prev_ideas_string}

Begin by generating an interestingly new high-level research proposal
that differs from what you have previously proposed.
...

# reflection prompt
Round {current_round}/{num_reflections}.

In your thoughts, first carefully consider the quality, novelty,
and feasibility of the proposal you just created.
Include any other factors that you think are important in evaluating
the proposal.
Ensure the proposal is clear and concise, and the JSON is in
```

- Each stage: launched with a specific human prompt
- Quite long prompt

# Parallelized Agentic Tree Search: Mimicking Human Exploration



- Integrates tree search with LLM-driven workflows [Chan et al.'25] across four expt stages
- Agentic tree search enables deeper and more systematic exploration of hypotheses

# Human Reviewer Feedback

**Reviewer #2: Compositional Regularization: Unexpected Obstacles in Enhancing Neural Network Generalization**

This paper investigates the effectiveness of incorporating a compositional regularization term into the loss function of neural networks to improve compositional generalization. The authors hypothesized that penalizing deviations from compositional structures would enhance the model's ability to generalize to unseen arithmetic expressions. However, their results on synthetic arithmetic datasets showed that compositional regularization did not lead to significant improvements and, in some cases, even hindered learning.

I think this paper greatly contributes to the workshops theme and fits into the scope. Moreover, it is a great example of challenges that occur during such approaches and could be interesting to discuss in the workshop setting. While I think that the authors should further broaden the experiments to other tasks in order to increase the generalizability of the findings, I would still recommend to accept the paper.

Rating: 6: Marginally above acceptance threshold
Award: No Award
Confidence: 2: The reviewer is willing to defend the evaluation, but it is quite likely that the reviewer did not understand central parts of the paper

# Outline

- Introduction

- A focus on AI Scientist-v2 – main principles

- **Investigating the next stages (foresight)**

- Conclusion

# Next AI "Technical" Stages

- Evolving very quickly, starting from LLMs (which are relatively recent)

- All research domains are gradually impacted
  - Depends on the quality of the Metropolis-Hasting-like proposal produced by the retrained LLM within the complex space related to the research domain under consideration (ex.: Coding, Mathematics)

- Quantity/quality of data for learning, computing resources

To investigate AlphaEvolve's breadth, we applied the system to over 50 open problems in mathematical analysis, geometry, combinatorics and number theory. The system's flexibility enabled us to set up most experiments in a matter of hours. In roughly 75% of cases, it rediscovered state-of-the-art solutions, to the best of our knowledge.

And in 20% of cases, AlphaEvolve improved the previously best known solutions, making progress on the corresponding open problems. For example, it advanced the kissing number problem. This geometric challenge has fascinated mathematicians for over 300 years and concerns the maximum number of non-overlapping spheres that touch a common unit sphere. AlphaEvolve discovered a configuration of 593 outer spheres and established a new lower bound in 11 dimensions.

AlphaEvolve (Google DeepMind)
May 14, 2025

# Positioning of the Future Human Scientist

- AI/LLM tools used frequently

Nombre de réponses obtenues : 318.

**1/ Quelles plateformes d'IA générative utilisez-vous dans le cadre professionnel ?**

Sur l'ensemble des réponses, ChatGPT dans ses différentes versions, arrive largement en tête (63%) ; il est utilisé seul ou avec Copilot (8%), DeepL (8%), Mistral (4%), Claude (2%). De façon plus anecdotique on a Le Chat (1%), LLama (1%), puis Grammarly, Perplexity, BlackBoxAI, Antidote ...

**2/ Pour quels usages et avec quel niveau de satisfaction ?**

L'usage principal des outils est 1) l'aide à la rédaction, correction grammaticale, synthèse de texte (29%), et 2) à la génération de code, explication de bugs, revue de code et apprentissage d'un langage informatique (28%) ; vient ensuite 3) la traduction( FR-UK, UK-FR) (22%).

Ces outils servent aussi à la recherche d'informations, à la documentation sur des sujets nouveaux (9%). Ils permettent de générer de nouvelles idées (4%). Ils sont aussi utilisés par curiosité pour explorer leurs capacités, les tester (4%).

- Will replace (at least) incremental research discoveries => focus on more complex problems
- Good AI-based research will rely on good human-based prompts => prompt engineering

# Regulations, Environmental and Other Considerations

- Classic questions related to the use of LLMs: data propriety, environmental impact, computing resources

- Authorship of AI-based research?

- High speed evolution of AI regulation (AI-act)

- Many research institutions publish their own internal recommendations

**RÉPUBLIQUE FRANÇAISE**
*Liberté*
*Égalité*
*Fraternité*

*Inria*

Note de cadrage sur l'usage de l'IA Générative au sein d'Inria

**Référence Gédéi : 17929**

**Date :** 12 mai 2025

**Rédacteur :** Jean-Frédéric Gerbeau

**Contributeurs :** CORUM[1]: Sylvain Petitjean ; DAFP : Catherine Gallet-Rybak, Elodie Limou, Alexandre Sicard ; DAJ : Eric Jaouen ; DGDI : David Rey ; DGDS : Karteek Alahari, Christophe Biernacki, Jean-Frédéric Gerbeau ; DPD : Anne Combe, Emilie Masson ; DSI : Florian Dufour, Thierry Murgue ; FSD : Didier Benza ; RSSI : Dominique Launay.

**Objet :** Consignes en matière d'usage d'outils d'IA générative

# Outline

- Introduction

- A focus on AI Scientist-v2 – main principles

- Investigating the next stages (foresight)

- **Conclusion**

- Impressive advances in AI-based research

- Be aware of these advances to
  - use these tools for some kinds of research (incremental research)
  - do human research otherwise (avoid incremental research)

- But these are just a few early ideas to be discussed. . .

# Thank you!