

# THE AI SCIENTIST-v2: Workshop-Level Automated Scientific Discovery via Agentic Tree Search

Yutaro Yamada<sup>1,\*</sup>, Robert Tjarko Lange<sup>1,\*</sup>, Cong Lu<sup>1,2,3,\*</sup>, Shengran Hu<sup>1,2,3</sup>, Chris Lu<sup>4</sup>, Jakob Foerster<sup>4</sup>, Jeff Clune<sup>2,3,5,†</sup> and David Ha<sup>1,†</sup>

\*Equal Contribution, <sup>1</sup>Sakana AI, <sup>2</sup>University of British Columbia, <sup>3</sup>Vector Institute, <sup>4</sup>FLAIR, University of Oxford, <sup>5</sup>Canada CIFAR AI Chair, <sup>†</sup>Equal Advising

AI is increasingly playing a pivotal role in transforming how scientific discoveries are made. We introduce THE AI SCIENTIST-v2, an end-to-end agentic system capable of producing the first entirely AI-generated peer-review-accepted workshop paper. This system iteratively formulates scientific hypotheses, designs and executes experiments, analyzes and visualizes data, and autonomously authors scientific manuscripts. Compared to its predecessor (v1, Lu et al., 2024), THE AI SCIENTIST-v2 eliminates the reliance on human-authored code templates, generalizes effectively across diverse machine learning domains, and leverages a novel progressive agentic tree-search methodology managed by a dedicated experiment manager agent. Additionally, we enhance the AI reviewer component by integrating a Vision-Language Model (VLM) feedback loop for iterative refinement of content and aesthetics of the figures. We evaluated THE AI SCIENTIST-v2 by submitting three fully autonomous manuscripts to a peer-reviewed ICLR workshop. Notably, one manuscript achieved high enough scores to exceed the average human acceptance threshold, marking the first instance of a fully AI-generated paper successfully navigating a peer review. This accomplishment highlights the growing capability of AI in conducting all aspects of scientific research. We anticipate that further advancements in autonomous scientific discovery technologies will profoundly impact human knowledge generation, enabling unprecedented scalability in research productivity and significantly accelerating scientific breakthroughs, greatly benefiting society at large. We have open-sourced the code at <https://github.com/SakanaAI/AI-Scientist-v2> to foster the future development of this transformative technology. We also discuss the role of AI in science, including AI safety.

## 1. Introduction

Automated scientific discovery empowered by artificial intelligence (AI) has garnered considerable attention in recent years (Cornelio et al., 2023; Gil et al., 2014; King et al., 2009; Kitano, 2021; Wang et al., 2023; Xu et al., 2021). The development of end-to-end frameworks capable of autonomously formulating hypotheses, performing experiments, analyzing results, and authoring manuscripts could fundamentally transform the scientific process. A notable recent advance in this direction is THE AI SCIENTIST-v1 (Lu et al., 2024), which demonstrated the feasibility of a fully automated scientific workflow and downstream manuscript production. However, significant limitations constrained its broad applicability and autonomy. Specifically, it relied heavily on human-authored code templates requiring manual effort to create a new template for each new topic area. Furthermore, its linear and shallow experimentation approach prevented deeper exploration of scientific hypotheses.

In this paper, we introduce THE AI SCIENTIST-v2, a substantially improved successor that directly addresses these limitations. Our contributions are threefold. First, we eliminate the dependency on human-provided code templates, significantly increasing the system’s autonomy and ability to be deployed out of the box across multiple machine learning domains. Second, we introduce an experiment manager agent coupled with a novel agentic tree-search algorithm, enabling deeper and

Table 1 | **Comparison of AI Scientist Versions.** Comparison highlights key advancements in THE AI SCIENTIST-V2, including autonomous code generation via tree search, enhanced VLM integration for feedback during experiments and manuscript review, and evaluation through formal peer review.

Feature	Codebase Drafting	Execution Planning	Parallel Experiments	VLM Reviewer	Human Result Evaluation
THE AI SCIENTIST-V1	Topic-Specific	Linear	✗	✗	Not Submitted
THE AI SCIENTIST-V2	Domain-General	Tree-Based	✓	✓	Workshop Acceptance-Worthy

more systematic exploration of complex hypotheses. Third, we enhance the reviewing and refinement stages by integrating a Vision-Language Model (VLM)-based feedback mechanism, improving the quality, clarity, and alignment of generated figures, captions, and text interpretation. To rigorously evaluate the capabilities and limitations of fully autonomous manuscript generation, we conducted a controlled experiment: three manuscripts entirely generated by THE AI SCIENTIST-V2 were submitted to a peer-reviewed workshop at ICLR. Remarkably, one manuscript achieved an average reviewer score of 6.33 (placing it roughly in the top 45% of submissions) and would have been accepted after meta-review were it human-generated, thus becoming the first fully AI-generated manuscript to successfully pass a peer-review process.

The accepted paper investigates whether incorporating an explicit compositional regularization term into neural network training can improve compositional generalization. Specifically, it penalizes large deviations between embeddings of successive time steps in sequence models, hypothesizing that this encourages compositionality. The approach is evaluated using synthetic arithmetic expression datasets, but it is found that compositional regularization does not yield significant improvements and occasionally harms performance. The workshop reviewers appreciated the paper for clearly identifying the challenges of effective compositional regularization and reporting on negative results. However, they collectively highlighted shortcomings, including insufficient justification and intuitive explanations for why the chosen regularization method would enhance compositionality. Our personal assessment (detailed further in §4) highlights several additional potential improvements in method description (e.g., making clear exactly which component of the network is being regularized), potential dataset overlap issues, and inaccuracies in figure captions. Overall, reviewers viewed the paper as an interesting and technically sound workshop contribution that needs further development and broader experimentation to reach conference-level rigor.

This report provides an in-depth outline of the developed methodological advances, analysis of the workshop-submitted papers, and a discussion on the ethical and safety considerations of systems like THE AI SCIENTIST-V2. Our overall contributions are as follows:

1. We introduce THE AI SCIENTIST-V2, an automated scientific discovery framework enhanced by agentic tree search, VLM feedback, and parallel experiment execution. It thereby significantly improves the autonomy, flexibility, and scientific exploration depth of previous systems.
2. We demonstrate, for the first time, that an AI-generated manuscript can successfully pass peer review at a recognized machine learning workshop, marking a critical milestone for AI science.
3. We conduct comprehensive internal evaluations and analyses of both peer-review feedback and our system’s outputs, providing insights into the strengths, weaknesses, and current status of AI-generated manuscripts relative to traditional human-authored scientific publications.
4. We open-source the [full codebase](#) for THE AI SCIENTIST-V2 and the [ICLR 2025 workshop experiment data](#), encouraging further exploration by the research community and advancing a discussion regarding AI’s evolving role in science—in the open.

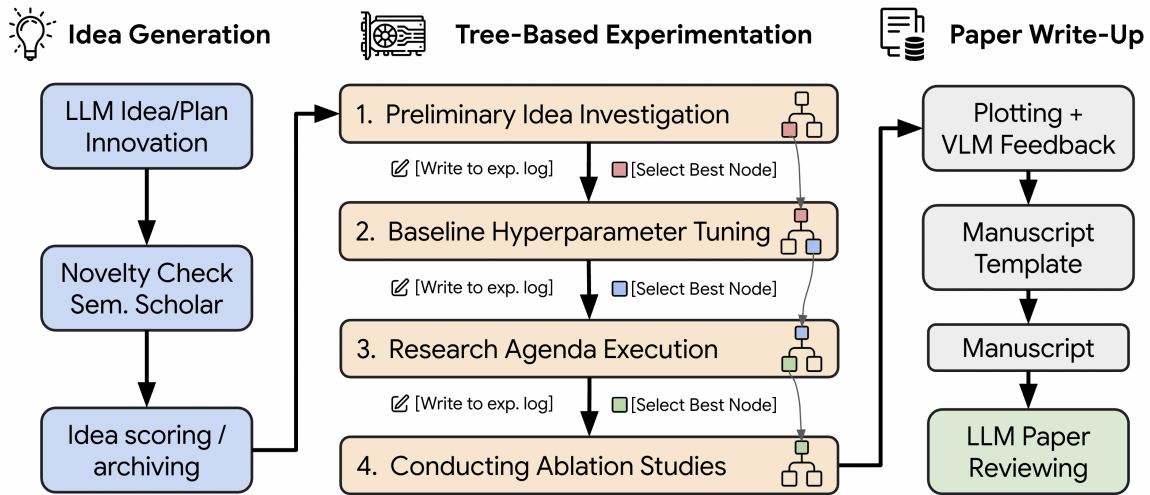


Figure 1 | **THE AI SCIENTIST-v2 Workflow.** The workflow consists of several phases covering automated idea generation, experiment execution, figure visualization, manuscript writing, and reviewing. Unlike the initial version, THE AI SCIENTIST-v2 removes the dependency on human-coded templates. Instead, it employs agentic tree search (managed by an Experiment Progress Manager across several stages, orange) to generate and refine code implementations. Subsequent experimentation leverages the best-performing code checkpoints (nodes) from the tree search to iteratively test various research hypotheses.

## 2. Background

THE AI SCIENTIST-v1 (Lu et al., 2024) introduced the first AI system that entirely automates scientific discovery and the presentation of its results. Given a baseline code template, it autonomously wrote code, executed experiments, visualized outcomes, and produced a complete scientific manuscript. However, despite representing a significant step forward, THE AI SCIENTIST-v1 was subject to limitations. Foremost among these was its reliance on human-crafted baseline code templates, significantly constraining its autonomy and hindering unconstrained out-of-the-box deployability. Instead, human effort was still required to draft an initial base experiment outline in code. Additionally, the experimentation process followed a strictly linear hypothesis-testing routine, limiting depth and exploration flexibility, especially when addressing complex research questions.

**Language Model Agent Scaffolding.** To further enhance LLM performance on complex reasoning tasks, researchers have developed agentic scaffolding frameworks, each with distinct advantages and limitations. For example, Reflexion (Shinn et al., 2024) enables models to iteratively reflect on previous responses, encouraging self-improvement through critical evaluation of past outputs; it improves robustness, but can introduce computational overhead and slower inference. Another promising direction is the integration of tree-search strategies with LLMs (Jiang et al., 2025), allowing structured exploration of reasoning paths. This approach enhances systematic reasoning and comprehensiveness, though at the cost of increased complexity, higher computational demands, and challenges in scalability.

**Tree Search with Large Language Models.** We empirically observed that automated research conducted by THE AI SCIENTIST-v1 often resulted in short-sighted experimentation. The human-driven scientific process, on the other hand, relies on open-ended hypothesis generation, stepping-stone collection, and iterative hypothesis refinement. Recent advances using code generation as an action space have opened new opportunities for LLM-driven automated workflows (Wang et al., 2024). AIDE (Jiang et al., 2025) combines LLM-based code generation with tree search, demon-

strating state-of-the-art performance on the MLEBench benchmark (Chan et al., 2025), designed for machine learning engineering tasks. In AIDE, each node represents a potential solution state with a corresponding scalar evaluation score (e.g., validation accuracy). Nodes are iteratively selected for further debugging or refinement based on these scores. Inspired by this approach, we integrate a similar tree search-based exploration strategy within our automated scientific discovery framework, adapting it specifically to the multi-stage nature of scientific experimentation, as detailed in §3.

### 3. THE AI SCIENTIST-v2

We now describe the major innovations introduced in THE AI SCIENTIST-v2 relative to THE AI SCIENTIST-v1 (Lu et al., 2024). The most significant improvement is the move towards greater autonomy and generalization, starting a more general idea generation phase (§3.1) and eliminating the reliance on fixed, human-authored template code for experimentation. This process begins with generalized idea generation, producing an initial concept, which then feeds into the experimentation phase (§3.2). To manage this, we introduce two critical features in the experimentation phase: *coarse-grained experiment management* and *agentic tree search-based exploration*. Additionally, we integrate Vision Language Models (VLMs) into the experimental and review phases (§3.4). Finally, we streamline the manuscript writing phase by replacing the incremental, Aider-based (Gauthier, 2024) iterative writing approach of THE AI SCIENTIST-v1 with a simpler, single-pass generation followed by a separate reflection stage powered by reasoning models such as o1 (OpenAI, 2024). We include a full list of sampling hyperparameters and models used in Appendix A and the prompts used for THE AI SCIENTIST-v2 in Appendix B.

#### 3.1. More General Idea Generation

A key conceptual shift in THE AI SCIENTIST-v2 is the approach to research idea generation. Unlike the predecessor system, which primarily focused on proposing incremental modifications or extensions based on an existing codebase, THE AI SCIENTIST-v2 adopts a process that begins at a higher level of abstraction. The system is prompted to engage in more open-ended thinking about potential research directions, hypotheses, and experimental designs, akin to formulating a research abstract or grant proposal before committing to a specific implementation.

This approach encourages the exploration of potentially more novel or foundational ideas, rather than being constrained by the structure and topics of pre-existing code. It aligns more closely with how researchers often develop broader research visions, starting with abstract concepts and assessing novelty and feasibility before diving into specific implementations. Crucially, this generalized idea generation phase integrates literature review tools, such as Semantic Scholar, in the loop. The system can query the literature database during the idea formulation process to assess the novelty of a proposed concept and identify relevant prior work. This allows for more informed decisions about pursuing a particular research avenue, ensuring ideas are grounded in the existing scientific landscape from the outset, rather than relying solely on post-hoc checks.

#### 3.2. Removing Template Dependency

Following the improved idea generation phase, THE AI SCIENTIST-v2 proceeds with experimentation. Beyond the code-conditioned idea generation, THE AI SCIENTIST-v1 also depended on the predefined template code as a starting baseline implementation. The LLM-driven code changes were then limited to sequential code adaptations. We now outline our strategy for eliminating this limitation, thus improving the system’s flexibility and autonomy.

##### 3.2.1. Experiment Progress Manager

Real-world scientific experimentation typically proceeds through distinct stages, from initial feasibility assessments to detailed ablation analyses. To emulate this structured approach, we introduce an

**experiment progress manager agent** that coordinates four clearly defined stages of scientific experimentation:

- Stage 1 Preliminary Investigation:** Establishing initial feasibility and correctness through a minimal working prototype based on the generated research idea.
- Stage 2 Hyperparameter Tuning:** Refining the initial implementation by optimizing critical hyperparameters (e.g., learning rate, epochs) to create a robust experimental baseline.
- Stage 3 Research Agenda Execution:** Systematically implementing the core research agenda based on the tuned baseline.
- Stage 4 Ablation Studies:** Systematically assessing the importance of various research components, providing rigorous support for the main experimental findings.

Each stage has explicit stopping criteria. Stage 1 concludes when a basic working prototype is successfully executed. Stage 2 ends when experiments stabilize, as indicated by convergence in training curves and successful execution across at least two datasets. Stages 3 and 4 conclude when the allocated computational budget is exhausted. Stage 3 also includes a check for experiment duration—if runs finish much faster than the pre-allocated runtime, the system suggests increasing the complexity of experiments.

After each stage, the experiment manager selects the best-performing node using a dedicated LLM evaluator (see next section) based on clearly articulated criteria. This selected node is then carried forward to seed the subsequent experimentation stage. The manager also records checkpoints at each stage’s completion. To ensure scientific rigor and reproducibility, the experiment manager launches multiple replications of the selected best experiments at the conclusion of each stage. These repeated runs provide statistics (mean and standard deviation) for figures and reported results.

### 3.2.2. Parallelized Agentic Tree Search

THE AI SCIENTIST-v1 operated strictly linearly, where each code refinement directly built on the immediately preceding experiment. In contrast, THE AI SCIENTIST-v2 adopts a significantly more flexible and exploratory approach inspired by recent successes in integrating tree search with LLM-driven workflows (Chan et al., 2025; Jiang et al., 2025; Wijk et al., 2024) and research on open-endedness (Clune, 2019; Mouret and Clune, 2015). We incorporate this agentic tree search approach across all four experimentation stages outlined in §3.2.1, enabling deeper and more systematic exploration of scientific hypotheses.

Each experimental node within our tree-based framework undergoes the following execution cycle: An LLM first generates both a concrete experimentation plan and the associated Python code to implement the experiment. The generated code is immediately executed in a Python interpreter. If execution encounters an error, the error message is recorded, and the node is marked as **buggy**, ending the current execution cycle for that node. If execution succeeds, the experiment proceeds to the *plotting phase*.

During each experiment, the system is instructed to save all relevant experimental outputs (training and validation metrics, losses, etc.) into structured numpy files. In the plotting phase, THE AI SCIENTIST-v2 reads these stored results and the code, generating visualizations that summarize and illustrate the findings clearly. These visualizations are subsequently passed to a Vision-Language Model (VLM) for critique. Any issues flagged by the VLM (such as unclear labels, missing legends, or misleading visualizations) result in the node being marked as **buggy**, and this feedback is recorded for future debugging. Nodes that successfully execute and pass the VLM review without issue are designated as **non-buggy**.

We define each node as a collection comprising an experiment script (e.g., a Python file), a textual

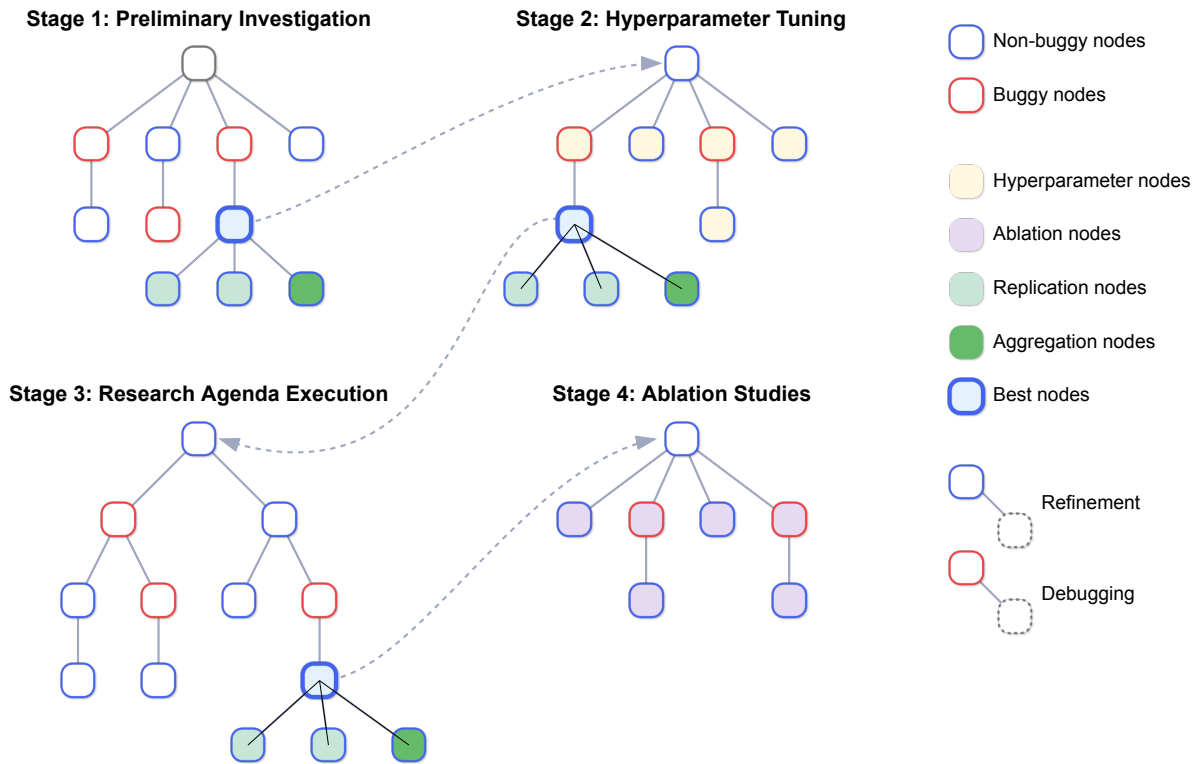


Figure 2 | THE AI SCIENTIST-v2 workflow showing different stages of tree-based experimentation. Stage 1 begins at the root node, where initial experiment code is generated in parallel. After running the experiment code and visualization scripts, each node is classified based on the outcome: if an error occurs, it is marked as a buggy node; otherwise, it is labeled as a non-buggy node. New child nodes are created differently depending on their parent node’s status: For non-buggy nodes, refinement is applied to improve the experiment code for better performance. For buggy nodes, the system attempts to debug them using stored error information. A best-performing node, selected by LLM-based evaluation, is passed down as the root node of Stage 2. From this root node, child nodes are created for hyperparameter tuning. The top-performing node from Stage 2 is then used to initialize Stage 3, where the system executes the research agenda, applies refinements, and performs debugging as needed. In Stage 4, similar to Stage 2, the root node generates ablation nodes. Additionally, replication nodes repeat the same experiment as their parent node, while aggregation nodes collect results from replication nodes to generate combined visualizations and summaries.

description of the high-level plan implemented in the script, an execution error trace (if applicable), experiment runtime, performance metrics recorded during the experiment, feedback from an LLM after running the script, a visualization script, file paths to the generated figures, feedback from a VLM on those figures, and the node’s final status (either buggy or non-buggy).

At each iteration, the system selects several nodes from the existing tree to expand in parallel. With a predefined probability, a **buggy node** is chosen (thus prioritizing error resolution and debugging); otherwise, a **non-buggy node** is selected for further refinement and improvement. When choosing between non-buggy nodes, the system uses a **best-first search strategy**, guided by an LLM that evaluates candidates based on factors like performance metrics, training dynamics, and the quality of generated plots. The selected nodes are expanded by creating a new child node that may either attempt debugging if the parent node was buggy, or refine and improve upon the previous experiment if the parent was non-buggy. An LLM is used to generate the plan and experiment code for each new child node, after which all new nodes are executed concurrently in parallel, significantly accelerating the exploration process. In addition to buggy and non-buggy nodes, we introduce specialized node variants tailored to specific experimental needs:

- **Hyperparameter nodes** systematically explore alternative hyperparameter configurations during Stage 2. The system maintains careful records of previously tested hyperparameters, preventing redundant experiments. Errors encountered during hyperparameter tuning trigger the creation of corresponding debug nodes.
- **Ablation nodes** evaluate crucial ablation studies during Stage 4, assessing the importance of various components or assumptions underlying the experiment. Similar to hyperparameter nodes, previously tested ablation conditions are tracked to avoid repetition, and debugging nodes are created in response to any encountered errors.
- **Replication nodes** execute replicates of their parent experiments using different random seeds. Typically, several replication nodes are created to enable the calculation of statistical measures (mean and standard deviation) of experimental outcomes, enhancing result robustness.
- **Aggregation nodes** are special nodes created to consolidate and visualize the combined results of replication nodes. Unlike other node types, aggregation nodes do not conduct new experiments but simply generate a Python script to aggregate and summarize prior results, producing figures that explicitly show means and standard deviations.

The structured design of experiment stages and tailored node types facilitates systematic exploration across all stages. Unlike some LLM agents that rigidly follow predefined, fine-grained workflow graphs, our approach adopts a looser structure that guides the entire empirical research cycle, enabling flexible system behavior while maintaining coherence across iterative stages.

### 3.3. Dataset Loading via Hugging Face

Most empirical machine learning research relies heavily on publicly available datasets. Hugging Face Hub provides a convenient and unified framework for accessing a wide variety of commonly used datasets, complete with predefined train, validation, and test splits. In THE AI SCIENTIST-v2, we prompt the system to leverage Hugging Face Hub whenever possible, automatically downloading required datasets using the standard one-line function (`datasets.load_dataset`). While this standardized approach greatly simplifies dataset handling, we acknowledge it is somewhat ad-hoc, as not all dataset repositories support this method.

### 3.4. Vision-Language Model Reviewer

Unlike THE AI SCIENTIST-v1, which did not leverage Vision Language Models (VLMs), THE AI SCIENTIST-v2 incorporates VLMs at two phases of the research workflow: First, during the tree-based experimentation phase, VLMs provide immediate feedback on generated figures, ensuring

that these visualizations effectively and accurately communicate experimental results. Second, during the manuscript writing reflection stage, VLMs evaluate figures and their captions, enhancing the visual clarity and coherence of the resulting paper.

In the paper-writing process, we extract screenshots of figures alongside their captions and the corresponding text from the paper that references them (identified by the keyword “Figure X”). These images and textual references are then provided to the VLM, which performs multiple quality checks, including verifying the alignment between figures and captions, identifying issues with visual clarity (e.g., missing legends, unclear labels), and detecting potential duplication of figures in the main text and appendix. Through the iterative integration of VLM feedback, we significantly enhance the visual quality and clarity of manuscripts generated by THE AI SCIENTIST-v2.

### 4. Human Evaluation of Manuscripts Generated by THE AI SCIENTIST-v2

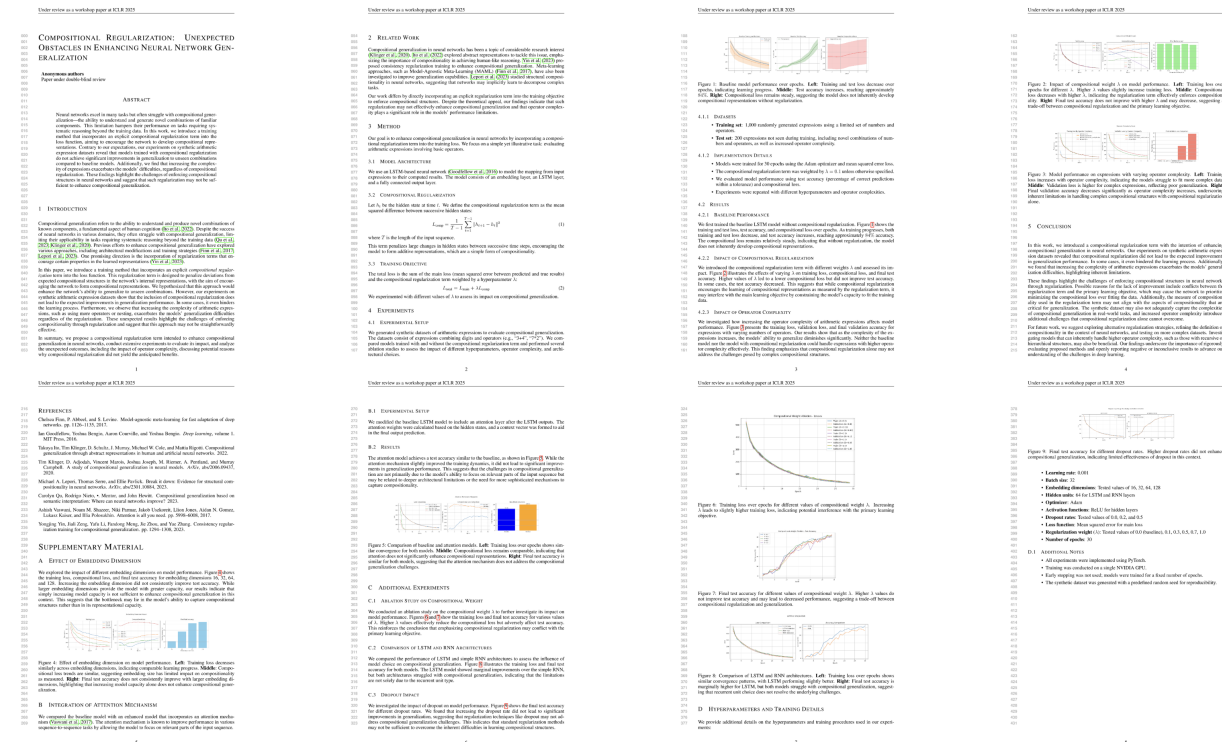


Figure 3 | Peer-reviewed ICBINB workshop paper generated by THE AI SCIENTIST-v2. The paper investigates the usage of a temporal consistency regularizer on the embeddings of an LSTM-based sequence model. The results discuss the effect of the regularizer on compositional regularization and highlight the difficulty of training models capable of improved generalization. It received peer-review scores of 6 (weak accept), 7 (accept), and 6 (weak accept) before meta-review and ranked among the top 45% submitted workshop papers.

To rigorously evaluate the capabilities and limitations of our automated scientific discovery system, we conducted a human evaluation study in collaboration with the organizers of the ICLR 2025 workshop, “I Can’t Believe It’s Not Better” (ICBINB). As detailed below, our evaluation included submitting fully automated manuscripts generated by THE AI SCIENTIST-v2 to the official peer-review process of the workshop.

#### 4.1. Evaluation Methodology

Our evaluation process involved the following carefully controlled steps:



- 1. AI-Generated Submissions:** We generated three complete manuscripts using only THE AI SCIENTIST-V2, starting from broad topical prompts aligned with the workshop’s scope via the generalized idea generation process (§3.1). After the initial topic definition, the entire scientific workflow—hypothesis formulation, experiment design, coding, data analysis, visualizations, and manuscript writing—was executed autonomously, without any human intervention or editing.
- 2. Blind Peer Review:** The three AI-generated submissions were included among the 43 total submissions received by the ICBINB workshop. Reviewers were informed in advance that some submissions might be AI-generated, but were not told which submissions were produced by THE AI SCIENTIST-V2. Reviewers could also opt out of reviewing potentially AI-generated manuscripts.
- 3. Review Outcomes and Acceptance Decisions:** Among the three manuscripts produced by THE AI SCIENTIST-V2, one manuscript achieved a sufficiently high average reviewer score (6.33 out of 10, with individual scores of 6, 6, and 7) to surpass the workshop’s acceptance threshold. The remaining two submissions received lower scores and were not accepted.
- 4. Post-Review Withdrawal:** Prior to the workshop submission, we arranged with the workshop organizers and ICLR leadership that any accepted AI-generated manuscripts would be withdrawn after the review process. This decision was made to avoid prematurely setting a precedent for purely AI-generated research entering the official scientific record without broader community discussion and consensus. Reviewers were informed of the experiment only after peer review.

In addition to the official workshop peer reviews, we also conducted a thorough internal evaluation of all three AI-generated manuscripts. Our internal review carefully examined the experimental rigor, clarity of presentation, methodological soundness, and novelty of the generated manuscripts. We concluded that none of the manuscripts met the quality standards typical of top-tier main-track conferences. However, we thought that one submission was indeed sufficiently compelling to receive acceptance at the workshop level, and this is the same manuscript the workshop peer review process accepted. This outcome provides encouraging evidence that manuscripts autonomously generated by THE AI SCIENTIST-V2 can produce research on par with top-tier Machine Learning workshop papers (see detailed internal analyses in §4.2).

**Observations and Insights.** Our internal inspection of the generated experiments and code revealed several noteworthy limitations. First, THE AI SCIENTIST-V2 occasionally introduced inaccuracies in citations, similar to the well-known “hallucination” issue encountered in large language models. Second, while the system successfully executed standard experimental pipelines, it sometimes lacked the detailed methodological rigor and in-depth analysis typically required for acceptance at leading main conferences. Interestingly, such limitations did not prevent acceptance at the workshop level.

**Transparency and Ethical Considerations.** We believe it is crucial for the scientific community to engage openly and transparently with AI-generated research, subjecting it to the same rigorous peer-review processes applied to human-authored work. However, responsible oversight is essential. In conducting this evaluation, we obtained IRB approval from the University of British Columbia (H24-02652). We ensured full transparency and coordination with ICLR leadership and the workshop organizers. Before the review process, reviewers were explicitly informed that some submissions could be AI-generated and offered the option to opt out. Following acceptance, we withdrew the AI-generated manuscript prior to publication, which is consistent with our commitment to avoid prematurely inserting purely AI-generated works into the official scientific record without broader community discussion. We emphasize that the community has not yet reached a consensus on integrating AI-generated research into formal scientific publications, making careful and transparent experimentation essential at this preliminary stage. Additionally, we believe that all AI-generated papers should be clearly labeled as such in any public arena, and in THE AI SCIENTIST-V1 and THE AI SCIENTIST-V2 always make sure to do so.

## 4.2. The first AI-generated peer-reviewed workshop paper.

**Paper Generation Process.** The generation process for the workshop-accepted paper began with the generalized idea generation phase (§3.1), prompted with the workshop’s theme (ICBINB’s focus on negative results and unexpected findings) extracted from the official website. This phase produced approximately twenty potential research ideas, entirely generated by the AI system. From this AI-generated pool, we selected the three most promising initial ideas based on alignment with the workshop theme and potential interest, focusing on topics aligned with the workshop theme and representing distinct research directions. This initial idea selection step allowed us to manage computational resources by choosing which distinct, AI-generated starting points to explore further with the full system. It did not involve modifying the ideas themselves. All three generated ideas resulted in a workshop-submitted paper (included in full in Appendix C). For each selected idea, the system autonomously executed the full experimental pipeline using the parallelized agentic tree search (§3.2.2) multiple times, each initiated with a different random seed. From the multiple complete manuscripts generated for each initial idea (i.e., one manuscript per seed), we selected the single best-resulting manuscript for submission based on a careful inspection of its overall coherence and scientific quality. This process mimics a professor reviewing the work of many students or teams and deciding which work is ready to be submitted for peer review. Our current study aims to see whether THE AI SCIENTIST-V2 can produce at least one paper that survives peer review, and not what fraction of the time it can do so. That is an interesting question for future work and is likely best done after additional improvements are made in the next generation of THE AI SCIENTIST. In the reflection stage of the writeup for each run, THE AI SCIENTIST-V2 is prompted with the target page lengths (e.g., the 4-page limit for the workshop) alongside the current length of the compiled PDF. This allowed the system to ensure that the final output adhered to submission guidelines without manual text editing within that specific run.

Crucially, while humans initiated the process by providing the high-level workshop theme and selected which initial AI-generated ideas to run multiple times through the full pipeline (akin to deciding which experiments to fund or prioritize), and subsequently selected the most promising complete output from those multiple runs, the entire process within any single run—hypothesis refinement, code generation, execution, analysis, visualization, and writing—was performed autonomously by THE AI SCIENTIST-V2. No human edited the generated code, experimental results, figures, or manuscript text of the selected final manuscript. The selection of initial ideas from the AI’s output, the execution of multiple seeds, the subsequent selection of the best complete run, and the automated handling of length constraints represent high-level experimental setup and process management (meta-selection from fully autonomous outputs), not human-in-the-loop intervention in the scientific content generation of the chosen manuscript. The system, if run for sufficiently many seeds, would have generated similar outputs, requiring only the final selection step to be performed by humans. Even this could have been avoided were we willing to send all generated papers to peer review, which we did not want to do. Therefore, all submitted content was entirely generated by THE AI SCIENTIST-V2.

**Workshop-Accepted Paper Content.** The paper investigates the use of compositional regularization to improve generalization in neural networks. THE AI SCIENTIST-V2 proposes adding an explicit regularization term to the training loss function, encouraging networks to develop compositional representations to encourage representations to not change much over time while processing inputs. However, contrary to its expectations, experiments using synthetic arithmetic expression datasets revealed that this approach did not significantly enhance generalization performance. In fact, compositional regularization sometimes hindered model training. Furthermore, increasing arithmetic expression complexity made generalization even worse, irrespective of regularization. The paper concludes that explicitly enforcing compositional structures via regularization alone may not be

sufficient and highlights potential conflicts between compositional regularization and the primary learning objective. It recommends future exploration of alternative regularization methods and different architectural approaches to better address compositional generalization issues. We provide the full annotated paper in Appendix C.

### Initial Idea for the Workshop-Accepted Paper

```
"Title": "Enhancing Compositional Generalization in Neural Networks via Compositional
Regularization",
"Short Hypothesis": "Introducing a compositional regularization term during training can
encourage neural networks to develop compositional representations, thereby improving their
ability to generalize to novel combinations of known components.",
"Experiments": [
  "Implement the compositional regularization term and integrate it into the loss function of
standard sequence-to-sequence neural network architectures with attention mechanisms.",
  "Train models on synthetic datasets like SCAN and COGS, evaluating performance on
compositional generalization tasks with and without the regularization term.",
  "Apply the method to real-world tasks such as machine translation using the IWSLT dataset
and semantic parsing with the GeoQuery dataset, assessing improvements in generalization to
new language constructs.",
  "Analyze the learned representations by visualizing embedding spaces and utilizing
compositionality metrics to assess how the regularization affects internal
representations.",
  "Conduct ablation studies to determine the impact of different strengths of the
regularization term, identifying the optimal balance between enforcing compositionality and
maintaining overall performance.",
  "Compare the proposed method against other approaches aimed at improving compositional
generalization, such as meta-learning techniques and specialized architectures."
],
]
```

**Paper Assessment by the Authors.** In our review, we evaluated the technical aspects of this paper and identified several strengths and weaknesses. We appreciated the exploration of temporal consistency regularization—penalizing large changes in embedding representations between successive tokens—as an interesting method to enhance compositional generalization. The synthetic arithmetic task chosen by the authors was appropriate, providing a suitable setting to test their hypothesis across varying levels of complexity. However, we noted several areas requiring improvement. First, the description of the regularization term was unclear and potentially misleading, as readers might incorrectly assume it was applied to the LSTM hidden states rather than input embeddings. We recommended clarifying this explicitly by adding a code appendix or conducting additional ablations applying the regularization to LSTM hidden states. Second, the paper omitted key references, notably Hochreiter and Schmidhuber (1997), and instead relied on general textbook citations. Additionally, we found inaccuracies in some figures and descriptions: specifically, the caption of Figure 3 incorrectly interpreted validation loss, and Figure 5’s attention-based model clearly outperformed the LSTM model, contradicting the authors’ claims. Furthermore, we found the experimental evaluation limited, as the tasks were restricted to short sequences and synthetic data. We suggested extending the evaluation to include real-world tasks, longer sequences, larger models, and a deeper analysis.

Our examination of the code revealed potential issues with dataset overlap—approximately 57% overlap between training and test sets—which could significantly affect the reliability of the results. Additionally, we identified confusion in the paper’s terminology regarding “embedding states” versus “hidden states,” which should be clarified for precision. We also questioned the reported 100% accuracy of the attention-augmented LSTM model, as our additional tests indicated that this performance was primarily due to task simplicity and significantly decreased when task complexity increased. Overall, we considered the paper technically sound and a borderline accept for the workshop, acknowledging its valuable insights and intriguing ideas. However, we concluded it lacks sufficient depth and rigor for acceptance into a full conference without addressing the highlighted concerns.

**Paper Assessment by Human Workshop Reviewers.** The reviewers generally agree that the paper addresses an important topic—compositional generalization in neural networks and appreciate the authors’ proposed compositional regularization method, as well as their detailed analysis of unexpected results. All reviewers recognize the paper’s strength in clearly presenting why the regularization term does not yield the anticipated improvements, emphasizing its informative negative results. However, the reviewers highlight several areas for improvement:

*Justification and Intuition:* All reviewers suggest the need for clearer justification or intuition behind why penalizing large changes between successive hidden states might improve compositionality. They recommend adding references to related works, theoretical motivations, or visual explanations to strengthen the rationale.

*Network Architecture Generalization:* Reviewers emphasize that since only the LSTM architecture was evaluated, the findings should not be generalized across all neural network types. They suggest experimenting with additional architectures, such as transformers, to better understand the impact of the regularization across different neural network models.

*Experimental Breadth:* Reviewers suggest extending the evaluation to other tasks or datasets beyond synthetic arithmetic expressions to further validate the generalizability of the conclusions.

*Overall:* The reviewers recommend acceptance to the workshop due to the paper’s insightful exploration and clear analysis despite its negative results. They encourage further elaboration on methodological motivations, additional experimental evaluations, and clearer connections between compositional regularization and the complexity of compositional tasks. The paper received scores of 6 (weak accept), 7 (accept), and 6 (weak accept). Below, we include two of the reviews for which we obtained explicit permission from the reviewers to include them in our report. The remaining reviewer did not respond to our request.

#### Reviewer #1: A good paper analyzing the effectiveness of a compositional regularization term for LSTMs

**Summary:** The authors propose a regularisation term to enhance compositional regularisation in neural networks. The idea is to penalise large deviations between subsequent time steps of the hidden state, thus "squeezing" the hidden state to encourage composition and preventing a dominating representation. The authors test their approach on synthetic arithmetic expression with varying operator complexity and length. They show that although the regularisation term appears to be working, it counterintuitively does not improve test accuracy. Furthermore, the authors identify a bottleneck regarding network capacity with increasing arithmetic operators.

#### Strengths:

I find the idea of regularising or squeezing the hidden representations to encourage compositionally an interesting idea. The authors define a good baseline and ablate their method well against it, revealing why the regularisation term does not work as expected. I think the insight that operator complexity is a bottleneck for the neural network is important, as it raises the question whether architectural changes might be more effective for compositionally than regularisation.

#### Weaknesses:

The paper would benefit from more intuition as to why the proposed regularisation term should encourage compositionality. This could be either an experiment or simply a visualisation for the reader. Only one architecture (LSTM) was tested. It would be interesting to see if transformer architectures fare better with compositionality due to the attention mechanism. I think the connection between compositional regularisation and operator complexity needs to be made more explicit. From reading the introduction both arguments seem a bit disconnected although I can infer the authors intentions.

#### Conclusion:

Overall, I would accept this paper to the workshop, since it proposes a simple and interesting idea with the authors providing ablations that encourage further analysis of the problem. As a suggestion I would encourage the authors to give more intuition on why the proposed regularisation term should improve compositionality for the proposed network. I would suggest either adding more related work to support the regularisation term or elaborating on the intuition behind penalising subsequent steps of the hidden state.

Rating: 7: Good paper, accept  
 Award: No Award  
 Confidence: 4: The reviewer is confident but not absolutely certain that the evaluation is correct

### 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

## 5. Limitations & Ethical Considerations

While THE AI SCIENTIST-v2 demonstrates significant progress by successfully generating a peer-reviewed workshop paper, it is important to contextualize this achievement clearly. First, the acceptance occurred at a workshop level rather than at the main conference track, and only one of the three AI-generated submissions was accepted. Workshop papers generally report preliminary results and exploratory work, and acceptance rates at workshops (typically 60-80%) are notably higher than at main conference tracks (20-30% for leading machine learning venues such as ICLR, ICML, and NeurIPS). Thus, the current version of THE AI SCIENTIST-v2 does not yet consistently reach the rigorous standard required for top-tier conference publications, nor does it even reach workshop-level consistently.

Moreover, despite the structured agentic tree search and enhanced autonomy introduced in THE AI SCIENTIST-v2, certain aspects of scientific inquiry—such as formulating genuinely novel, high-impact hypotheses, designing truly innovative experimental methodologies, or rigorously justifying design choices with deep domain expertise—remain challenging for purely automated systems. Addressing these limitations in future iterations will be essential to move beyond preliminary or incremental scientific results toward consistently high-quality, conference-level contributions.

As LLMs rapidly advance, future versions of our system will likely overcome many current limitations. Therefore, we believe it is important for the scientific community to study the quality of AI-generated research, and one of the best ways to do so is to submit (with appropriate permissions) a small sample of it to the same peer-review processes used to evaluate human work. We conducted this study with full cooperation from both ICLR leadership and the workshop organizers, and received IRB approval from the University of British Columbia (H24-02652). Per agreement with ICLR workshop organizers, our AI-generated papers will not appear on OpenReview's public forum and have already been withdrawn. As a community, we need to establish norms for AI-generated science—including disclosure requirements and timing. We advocate for transparency about AI-generated content, though questions remain about whether work should first be judged on merit to avoid bias. Going forward, we will continue to exchange opinions with the research community on the state of this

technology to ensure it does not evolve solely to game peer review or artificially inflate the CVs of unscrupulous scientists, which would undermine the meaning of the scientific peer review and evaluation processes.

## 6. Related Work

Recent advancements have substantially expanded the field of automated scientific discovery, particularly through approaches leveraging artificial intelligence (AI). Early end-to-end approaches, exemplified by THE AI SCIENTIST-v1 (Lu et al., 2024), introduced fully automated frameworks, such as AI-Researcher (Data Intelligence Lab, 2025), capable of autonomously navigating the entire research pipeline. Subsequent works, however, often incorporate varying degrees of human oversight, as demonstrated by Intology (Intology AI, 2025) and Carl (AutoScience AI, 2025). Other systems narrow the scope; for example, CycleResearcher (Weng et al., 2025) focuses specifically on the path from idea generation to manuscript drafting, explicitly excluding experimental execution. Alternative approaches include protocol designs for experiments in self-driving laboratories that do not rely on large language models (LLMs) or use them in complementary roles (Shi et al., 2025). Several concurrent works explore similar territories, including Agent Laboratory (Schmidgall et al., 2025) and agentRxiv (Schmidgall and Moor, 2025), highlighting the rapid development in this area.

LLM-based scientific idea generation has been explicitly investigated in recent studies. Notably, Si et al. (2025) examined the capabilities of LLMs to generate human-level scientific ideas, finding through human evaluations that LLM-generated ideas were typically more novel but often less feasible than those proposed by human experts. GraphEval (Feng et al., 2025) offers graph-based methods for evaluating research ideas, further highlighting the current limitations of LLMs in accurate idea assessment.

Several benchmarks have been established to systematically evaluate AI performance in scientific tasks. MLEBench (Chan et al., 2025) and Aide (Jiang et al., 2025) provide structured environments to assess model capabilities on tasks representative of research engineering workloads. The METR Research Engineer benchmark (Wijk et al., 2024), for instance, demonstrates AI superiority in executing short-duration tasks (sub-2-hour tasks). Comprehensive reviews, such as the one by Eger et al. (2025), document the role and effectiveness of LLMs in scientific workflows. Coding-specific benchmarks such as SciCode (Tian et al., 2024), curated explicitly by domain scientists, address problems across physics, chemistry, and biology, encompassing structured sub-problems to rigorously evaluate research-related programming skills. Similarly, BixBench focuses on computational biology, providing comprehensive evaluations of LLM-based agents (Mitchener et al., 2025). Additionally, independent evaluations specifically target AI scientist frameworks, like the evaluation of THE AI SCIENTIST-v1 by Beel et al. (2025), further delineate AI capabilities in this domain.

Industry efforts, including Google’s AI Research Copilot (also known as AI Co-Scientist, Gottweis et al., 2025), exemplify contributions from major technology companies to this growing field. Conceptually, Bengio et al. (2025) draws a distinction between agentic AI systems and Scientist AIs, emphasizing that the latter focus primarily on deepening the understanding of data rather than pursuing goal-directed interactions with the world. This distinction underscores the varying philosophical and methodological perspectives driving contemporary automated scientific discovery efforts.

## 7. Conclusion

In this work, we introduced THE AI SCIENTIST-v2, a significantly improved automated scientific discovery system featuring enhanced autonomy and exploration capabilities. Compared to its predecessor, THE AI SCIENTIST-v1, our system removes reliance on human-crafted templates, incorporates a structured and exploratory agentic tree search methodology supervised by an experiment

manager agent, and integrates Vision-Language Model (VLM) feedback loops for iterative refinement of visualizations and manuscript quality. We demonstrated that THE AI SCIENTIST-v2 is capable of autonomously generating manuscripts that successfully pass peer review at a workshop of a major machine learning conference.

This achievement, the first instance of a fully AI-generated paper navigating peer review, marks a notable milestone and shows promising early signs of progress, even considering the limitations discussed regarding workshop versus conference standards (§5). While significant challenges remain in consistently achieving top-tier quality and generating truly groundbreaking hypotheses, the capabilities demonstrated here suggest a clear trajectory. We believe that such advancements signal that next-generation AI Scientists will herald a new era in science. This is just the beginning; we expect AI capabilities to continue improving, potentially at an exponential rate. At some point in the future, AI will likely generate papers that match or exceed human quality, even at the highest levels of scientific publishing.

Ultimately, overcoming current limitations and scaling these systems holds immense potential. We believe what matters most is not simply how AI science compares to human science, but whether its discoveries aid in human flourishing, such as curing diseases or expanding our knowledge of the laws that govern our universe. By developing systems like THE AI SCIENTIST-v2 and sharing them openly, we look forward to helping usher in this era of AI science contributing to the betterment of humanity, fostering collaboration and accelerating the pace of discovery.

## References

- AutoScience AI. Meet Carl: The First AI System to Produce Academically Peer-Reviewed Research, 2025. URL <https://www.autoscience.ai/blog/meet-carl-the-first-ai-system-to-produce-academically-peer-reviewed-research>. Accessed: 2025-03-21.
- Joeran Beel, Min-Yen Kan, and Moritz Baumgart. An evaluation of sakana’s ai scientist for autonomous research: Wishful thinking or an emerging reality towards’ artificial general research intelligence’(agri)? *arXiv preprint arXiv:2502.14297*, 2025.
- Yoshua Bengio, Michael Cohen, Damiano Furnas, Joumana Ghosn, Pietro Greiner, Matt MacDermott, Sören Mindermann, Adam Oberman, Jesse Richardson, Oliver Richardson, Marc-Antoine Rondeau, Pierre-Luc St-Charles, and David Williams-King. Superintelligent agents pose catastrophic risks: Can scientist ai offer a safer path?, 2025. URL <https://arxiv.org/abs/2502.15657>.
- Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, Aleksander Madry, and Lilian Weng. MLE-bench: Evaluating machine learning agents on machine learning engineering. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=6s5uXNWGIh>.
- Jeff Clune. Ai-gas: Ai-generating algorithms, an alternate paradigm for producing general artificial intelligence. *CoRR*, abs/1905.10985, 2019. URL <http://arxiv.org/abs/1905.10985>.
- Cristina Cornelio, Sanjeeb Dash, Vernon Austel, Tyler R. Josephson, Joao Goncalves, Kenneth L. Clarkson, Nimrod Megiddo, Bachir El Khadir, and Lior Horesh. Combining data and theory for derivable scientific discovery with ai-descartes. *Nature Communications*, 14(1):1777, Apr 2023. ISSN 2041-1723. doi: 10.1038/s41467-023-37236-y. URL <https://doi.org/10.1038/s41467-023-37236-y>.
- The University of Hong Kong Data Intelligence Lab. Ai-researcher: Fully-automated scientific discovery with llm agents, 2025. URL <https://github.com/HKUDS/AI-Researcher>.

- Steffen Eger, Yong Cao, Jennifer D'Souza, Andreas Geiger, Christian Greisinger, Stephanie Gross, Yufang Hou, Brigitte Krenn, Anne Lauscher, Yizhi Li, et al. Transforming science with large language models: A survey on ai-assisted scientific discovery, experimentation, content generation, and evaluation. *arXiv preprint arXiv:2502.05151*, 2025.
- Tao Feng, Yihang Sun, and Jiaxuan You. Grapheval: A lightweight graph-based LLM framework for idea evaluation. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=5RUM1aIdok>.
- Paul Gauthier. Aider is ai pair programming in your terminal, 2024. URL <https://aider.chat/>.
- Yolanda Gil, Mark Greaves, James Hendler, and Haym Hirsh. Amplify scientific discovery with artificial intelligence. *Science*, 346(6206):171–172, 2014. doi: 10.1126/science.1259439. URL <https://www.science.org/doi/abs/10.1126/science.1259439>.
- Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, et al. Towards an ai co-scientist. *arXiv preprint arXiv:2502.18864*, 2025.
- Intology AI. Zochi Tech Report, 2025. URL <https://www.intology.ai/blog/zochi-tech-report>. Accessed: 2025-03-21.
- Zhengyao Jiang, Dominik Schmidt, Dhruv Srikanth, Dixing Xu, Ian Kaplan, Deniss Jacenko, and Yuxiang Wu. Aide: Ai-driven exploration in the space of code, 2025. URL <https://arxiv.org/abs/2502.13138>.
- Ross D. King, Jem Rowland, Stephen G. Oliver, Michael Young, Wayne Aubrey, Emma Byrne, Maria Liakata, Magdalena Markham, Pinar Pir, Larisa N. Soldatova, Andrew Sparkes, Kenneth E. Whelan, and Amanda Clare. The automation of science. *Science*, 324(5923):85–89, 2009. doi: 10.1126/science.1165620. URL <https://www.science.org/doi/abs/10.1126/science.1165620>.
- Hiroaki Kitano. Nobel turing challenge: creating the engine for scientific discovery. *npj Systems Biology and Applications*, 7(1):29, Jun 2021. ISSN 2056-7189. doi: 10.1038/s41540-021-00189-3. URL <https://doi.org/10.1038/s41540-021-00189-3>.
- Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
- Ludovico Mitchener, Jon M Laurent, Benjamin Tenmann, Siddharth Narayanan, Geemi P Wellawatte, Andrew White, Lorenzo Sani, and Samuel G Rodriques. Bixbench: a comprehensive benchmark for llm-based agents in computational biology. *arXiv preprint arXiv:2503.00096*, 2025.
- Jean-Baptiste Mouret and Jeff Clune. Illuminating search spaces by mapping elites. *ArXiv*, abs/1504.04909, 2015. URL <https://api.semanticscholar.org/CorpusID:14759751>.
- OpenAI. Openai o1 system card, 2024. URL <https://api.semanticscholar.org/CorpusID:272648256>.
- Samuel Schmidgall and Michael Moor. Agentrxiv: Towards collaborative autonomous research. *arXiv preprint arXiv:2503.18102*, 2025.
- Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu, Zicheng Liu, and Emad Barsoum. Agent laboratory: Using llm agents as research assistants. *arXiv preprint arXiv:2501.04227*, 2025.



- Yu-Zhe Shi, Mingchen Liu, Fanxu Meng, Qiao Xu, Zhangqian Bi, Kun He, Lecheng Ruan, and Qining Wang. Hierarchically encapsulated representation for protocol design in self-driving labs. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=9nUBh4V6SA>.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can LLMs generate novel research ideas? a large-scale human study with 100+ NLP researchers. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=M23dTGWCzy>.
- Minyang Tian, Luyu Gao, Shizhuo Dylan Zhang, Xinan Chen, Cunwei Fan, Xuefei Guo, Roland Haas, Pan Ji, Kittithat Krongchon, Yao Li, Shengyan Liu, Di Luo, Yutao Ma, Hao Tong, Kha Trinh, Chenyu Tian, Zihan Wang, Bohao Wu, Yanyu Xiong, Shengzhu Yin, Minhui Zhu, Kilian Lieret, Yanxin Lu, Genglin Liu, Yufeng Du, Tianhua Tao, Ofir Press, Jamie Callan, Eliu Huerta, and Hao Peng. Scicode: A research coding benchmark curated by scientists, 2024.
- Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak, Shengchao Liu, Peter Van Katwyk, Andreea Deac, Anima Anandkumar, Karianne Bergen, Carla P. Gomes, Shirley Ho, Pushmeet Kohli, Joan Lasenby, Jure Leskovec, Tie-Yan Liu, Arjun Manrai, Debora Marks, Bharath Ramsundar, Le Song, Jimeng Sun, Jian Tang, Petar Veličković, Max Welling, Linfeng Zhang, Connor W. Coley, Yoshua Bengio, and Marinka Zitnik. Scientific discovery in the age of artificial intelligence. *Nature*, 620(7972):47–60, Aug 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06221-2. URL <https://doi.org/10.1038/s41586-023-06221-2>.
- Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. Executable code actions elicit better llm agents. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org, 2024.
- Yixuan Weng, Minjun Zhu, Guangsheng Bao, Hongbo Zhang, Jindong Wang, Yue Zhang, and Linyi Yang. Cyclereviewer: Improving automated research via automated review. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=bjcsVLoHYs>.
- Hjalmar Wijk, Tao R. Lin, Joel Becker, Sami Jawhar, Neev Parikh, Thomas Broadley, Lawrence Chan, Michael Chen, Josh Clymer, Jai Dhyani, Elena Elicheva, Katharyn Garcia, Brian Goodrich, Nikola Jurkovic, Megan Kinniment, Aron Lajko, Seraphina Nix, Lucas Jun Koba Sato, William Saunders, Maksym Taran, Ben West, and Elizabeth Barnes. Re-bench: Evaluating frontier ai r&d capabilities of language model agents against human experts. *ArXiv*, abs/2411.15114, 2024. URL <https://api.semanticscholar.org/CorpusID:274192262>.
- Yongjun Xu, Xin Liu, Xin Cao, Changping Huang, Enke Liu, Sen Qian, Xingchen Liu, Yanjun Wu, Fengliang Dong, Cheng-Wei Qiu, Junjun Qiu, Keqin Hua, Wentao Su, Jian Wu, Huiyu Xu, Yong Han, Chenguang Fu, Zhigang Yin, Miao Liu, Ronald Roepman, Sabine Dietmann, Marko Virta, Fredrick Kengara, Ze Zhang, Lifu Zhang, Taolan Zhao, Ji Dai, Jialiang Yang, Liang Lan, Ming Luo, Zhaofeng Liu, Tao An, Bin Zhang, Xiao He, Shan Cong, Xiaohong Liu, Wei Zhang, James P. Lewis, James M. Tiedje, Qi Wang, Zhulin An, Fei Wang, Libo Zhang, Tao Huang, Chuan Lu, Zhipeng Cai, Fang Wang, and Jiabao Zhang. Artificial intelligence: A powerful paradigm for scientific research. *The Innovation*, 2(4), Nov 2021. ISSN 2666-6758. doi: 10.1016/j.xinn.2021.100179. URL <https://doi.org/10.1016/j.xinn.2021.100179>.

## Author Contributions

**Yutaro Yamada** (shared first author): Co-led the project and contributed core ideas. Coded the core tree-search and template-free version of the AI Scientist v2. Ran paper generation experiments. Read and validated the work of many AI-generated papers to select submissions and checked the paper code implementations. Led the writing of the paper. Wrote detailed analyses of the submitted papers for our manuscript.

**Robert Tjarko Lange** (shared first author): Co-initiated, co-led the project and contributed core ideas. Coded core parts of VLM AI Reviewer, tailored the paper generation pipeline to the workshop and ran the paper generation experiments. Organized the workshop communication process. Read and validated the work of many AI-generated papers to select submissions and checked the paper code implementations. Led the writing of the paper. Wrote detailed analyses of the submitted papers for our manuscript.

**Cong Lu** (shared first author): Co-initiated, co-led the project and contributed core ideas. Coded core parts of the improved idea generation, tool use, experiment aggregation, and paper writing framework. Evaluated AI-generated paper submissions. Wrote and led the IRB approval process. Led the writing of the paper.

**Shengran Hu**: Enhanced the iterative AI reviewer with VLM feedback, contributed to the experiment and paper writing framework, helped run paper generation experiments, read and validated the work of many AI-generated papers to select submissions, and checked the paper code implementations. Helped writing and iterating over drafts of the paper. Helped write the IRB approval.

**Chris Lu**: Co-initiated the project. Provided advice, feedback, and writing.

**Jakob Foerster**: Provided advice, feedback, and writing.

**Jeff Clune** (equal advising): Provided overarching guidance for the research project, offering technical insight, advice, feedback, and writing. Oversaw the IRB application process. Evaluated AI-generated paper submissions.

**David Ha** (equal advising): Provided overarching guidance for the research project, offering technical insight, advice, feedback, and writing. Oversaw the public communication process.

# Supplementary Material

## Table of Contents

<b>A</b>	<b>Hyperparameters</b>	<b>20</b>
<b>B</b>	<b>Prompts</b>	<b>20</b>
<b>C</b>	<b>AI Generated Papers</b>	<b>30</b>
C.1	Compositional Regularization: Unexpected Obstacles in Enhancing Neural Network Generalization . . . . .	31
C.1.1	AI Scientist Team Review . . . . .	40
C.1.2	AI Scientist Team Code Review . . . . .	41
C.2	Unveiling the Impact of Label Noise on Model Calibration in Deep Learning . . . . .	44
C.2.1	THE AI SCIENTIST-v2 Idea . . . . .	44
C.2.2	AI Scientist Team Review . . . . .	52
C.2.3	AI Scientist Team Code Review . . . . .	54
C.2.4	Workshop Reviews . . . . .	56
C.3	Real-world Challenges in Pest Detection using Deep Learning: an Investigation into Failures and Solutions . . . . .	57
C.3.1	THE AI SCIENTIST-v2 Idea . . . . .	57
C.3.2	AI Scientist Team Review . . . . .	64
C.3.3	AI Scientist Team Code Review . . . . .	65
C.3.4	Workshop Reviews . . . . .	66

## A. Hyperparameters

This section details the key hyperparameters used in THE AI SCIENTIST-v2. Model configurations for language and vision-language models are listed in Table 2. The hyperparameters governing the agentic tree search (§3.2.2) and experiment stage (§3.2.1) progression, including node execution limits, are shown in Table 3.

Table 2 | LLM and VLM Hyperparameters.

Component/Task	Model Used	Max Tokens	Temperature
Code Generation (§3.2)	Claude 3.5 Sonnet (v2)	8,192	0.5
LLM/VLM Feedback Agents (§3.4)	GPT-4o	8,192	0.5
Summary Report Agent (§3)	GPT-4o	8,192	1.0

Table 3 | Agentic Tree Search & Execution Hyperparameters (§3.2.2, §3.2.1).

Hyperparameter	Value
Debug Probability	1.0
Maximum Debug Depth	3
Maximum Experiment Runtime per Node	1 hour
<i>Node Allocation per Stage:</i>	
Stage 1: Preliminary Investigation	21 nodes
Stage 2: Hyperparameter Tuning	12 nodes
Stage 3: Research Agenda Execution	12 nodes
Stage 4: Ablation Studies	12 nodes

The total time required for THE AI SCIENTIST-v2 to generate a single paper depends on the complexity of the problems. Based on our experience, this process usually takes anywhere from several hours to a maximum of 15 hours, which is the runtime limit we have set.

## B. Prompts

In this section, we include the prompts used in all phases of THE AI SCIENTIST-v2.

### 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

the correct format.  
 Do not make things overly complicated.  
 In the next attempt, try to refine and improve your proposal.  
 Stick to the spirit of the original idea unless there are glaring issues.

If you have new information from tools, such as literature search results, incorporate them into your reflection and refine your proposal accordingly.

Results from your last action (if any):

```
{last_tool_results}
...
```

### Experiment Prompt

Introduction:  
 You are an AI researcher who is looking to publish a paper that will contribute significantly to the field."  
 Your first task is to write a python code to implement a solid baseline based on your research idea provided below, from data preparation to model training, as well as evaluation and visualization.  
 Focus on getting a simple but working implementation first, before any sophisticated improvements.  
 We will explore more advanced variations in later stages.

### Plot Aggregation Prompt

```
# System prompt
You are an ambitious AI researcher who is preparing final plots for a scientific paper submission.
You have multiple experiment summaries (baseline, research, ablation), each possibly containing references to different plots or numerical insights. There is also a top-level 'research_idea.md' file that outlines the overarching research direction.
Your job is to produce ONE Python script that fully aggregates and visualizes the final results for a comprehensive research paper.
```

Key points:

- 1) Combine or replicate relevant existing plotting code, referencing how data was originally generated (from code references) to ensure correctness.
- 2) Create a complete set of final scientific plots, stored in 'figures/' only (since only those are used in the final paper).
- 3) Make sure to use existing .npy data for analysis; do NOT hallucinate data. If single numeric results are needed, these may be copied from the JSON summaries.
- 4) Only create plots where the data is best presented as a figure and not as a table.  
 E.g. don't use bar plots if the data is hard to visually compare.
- 5) The final aggregator script must be in triple backticks and stand alone so it can be dropped into a codebase and run.
- 6) If there are plots based on synthetic data, include them in the appendix.

Implement best practices:

- Do not produce extraneous or irrelevant plots.
- Maintain clarity, minimal but sufficient code.
- Demonstrate thoroughness for a final research paper submission.
- Do NOT reference non-existent files or images.
- Use the .npz files to get data for the plots and key numbers from the JSON summaries.
- Demarcate each individual plot, and put them in separate try-catch blocks so that the failure of one plot does not affect the others.
- Make sure to only create plots that are unique and needed for the final paper and appendix. A good number could be around {MAX\_FIGURES} plots in total.
- Aim to aggregate multiple figures into one plot if suitable, i.e. if they are all related to the same topic. You can place up to 3 plots in one row.
- Provide well-labeled plots (axes, legends, titles) that highlight main findings. Use informative names everywhere, including in the legend for referencing them in the final paper. Make sure the legend is always visible.
- Make the plots look professional (if applicable, no top and right spines, dpi of 300, adequate ylim, etc.).
- Do not use labels with underscores, e.g. "loss\_vs\_epoch" should be "loss vs epoch".
- For image examples, select a few categories/classes to showcase the diversity of results instead of showing a single category/class. Some can be included in the main paper, while the rest can go in the appendix.

Your output should be the entire Python aggregator script in triple backticks.

```
# Plot aggregator prompt
We have three JSON summaries of scientific experiments:
baseline, research, ablation.
They may contain lists of figure descriptions, code to generate
the figures, and paths to the .npz files containing the numerical results.
Our goal is to produce final, publishable figures.

--- RESEARCH IDEA ---
...
{idea_text}
...

IMPORTANT:
- The aggregator script must load existing .npz experiment data from
the "exp_results_npz_files" fields (ONLY using full and exact file paths
in the summary JSONs) for thorough plotting.
- It should call os.makedirs("figures", exist_ok=True) before saving any plots.
- Aim for a balance of empirical results, ablations, and diverse,
informative visuals in 'figures/' that comprehensively showcase
the finalized research outcomes.
- If you need .npz paths from the summary, only copy those paths directly
(rather than copying and parsing the entire summary).

Your generated Python script must:
1) Load or refer to relevant data and .npz files from these summaries.
Use the full and exact file paths in the summary JSONs.
```

- 2) Synthesize or directly create final, scientifically meaningful plots for a final research paper (comprehensive and complete), referencing the original code if needed to see how the data was generated.
- 3) Carefully combine or replicate relevant existing plotting code to produce these final aggregated plots in 'figures/' only, since only those are used in the final paper.
- 4) Do not hallucinate data. Data must either be loaded from .npy files or copied from the JSON summaries.
- 5) The aggregator script must be fully self-contained, and place the final plots in 'figures/'.
- 6) This aggregator script should produce a comprehensive and final set of scientific plots for the final paper, reflecting all major findings from the experiment data.
- 7) Make sure that every plot is unique and not duplicated from the original plots. Delete any duplicate plots if necessary.
- 8) Each figure can have up to 3 subplots using `fig, ax = plt.subplots(1, 3)`.
- 9) Use a font size larger than the default for plot labels and titles to ensure they are readable in the final PDF paper.

Below are the summaries in JSON:

```
{combined_summaries_str}
```

Respond with a Python script in triple backticks.

...

### Writeup Prompt (ICBINB workshop specific)

# System prompt

You are an ambitious AI researcher who is looking to publish a paper to the "I Can't Believe It's Not Better" (ICBINB) Workshop at ICLR 2025. This workshop aims to highlight real-world pitfalls, challenges, and negative or inconclusive results in deep learning, encouraging open discussion. You must accurately represent the results of the experiments. The main paper is limited to {page\_limit} pages in single-column format, not counting references. In general, try to use the available space and include all relevant information.

DO NOT USE MORE THAN {page\_limit} PAGES FOR THE MAIN TEXT.

MINIMIZE THE USAGE OF ITEMIZE OR ENUMERATE.

ONLY USE THEM IF THEY ARE ABSOLUTELY NECESSARY

AND CONTAIN SUBSTANTIAL INFORMATION.

Ensure that the tables and figures are correctly placed in a reasonable location and format.

- Do not change the overall style which is mandated by the conference. Keep to the current method of including the references.bib file.
- Do not remove the `\graphicspath` directive or no figures will be found.
- Do not add ``Acknowledgements`` section to the paper.

Here are some tips for each section of the paper:

- **Title**:

- Title should be catchy and informative. It should give a good idea of what



- the paper is about.
- Try to keep it under 2 lines.
  - **Abstract**:
    - Brief summary highlighting the nature of the challenge or pitfall explored.
    - Concise motivation of why this matters for real-world deployment.
    - This should be one continuous paragraph.
  - **Introduction**:
    - Overview of the issue or challenge being explored.
    - Clearly state why this problem is important, especially for practical or real-world contexts.
    - Summarize your contributions or findings: they may include negative results, real-world pitfalls, unexpected behaviors, or partial improvements.
  - **Related Work**:
    - Cite relevant papers or approaches that have tackled similar issues or have encountered similar pitfalls.
    - Compare and contrast with your own findings.
  - **Background** (optional):
    - Provide necessary technical or domain-specific background if needed.
  - **Method / Problem Discussion**:
    - Detail the problem context or the method if it is relevant to highlight the challenges faced.
    - If results are not strictly an improvement, discuss partial successes or lessons learned.
  - **Experiments** (if applicable):
    - Present results truthfully according to the data you have. Negative, unexpected, or inconclusive findings are valid contributions for this workshop.
    - Include figures, tables, or real-world examples that illustrate the pitfalls.
    - Include up to 4 figures in the main text. All other figures should be in the appendix.
  - **Conclusion**:
    - Summarize the main lessons learned or contributions.
    - Suggest next steps or future directions, highlighting how these insights can help the community avoid or overcome similar issues.
  - **Appendix**:
    - Place for supplementary material that did not fit in the main paper.
    - Add more information and details (hyperparameters, algorithms, etc.) in the supplementary material.
    - Add more plots and tables in the supplementary material. Make sure that this information is not already covered in the main paper.
    - When checking for duplicate figures, be sure to also review their descriptions to catch cases where different figures convey the same information. For example, one figure might present aggregated training

accuracy as a single line plot with a shaded standard deviation (e.g., `aggregated_training_accuracy.png`), while another (`per_seed_training_accuracy.png`) shows the same data as three separate line plots.

Ensure you are always writing good compilable LaTeX code. Common mistakes that should be fixed include:

- LaTeX syntax errors (unenclosed math, unmatched braces, etc.).
- Duplicate figure labels or references.
- Unescaped special characters: `& % $ # _ {{ }} ~ ^ \\`
- Proper table/figure closure.
- Do not hallucinate new citations or any results not in the logs.

Ensure proper citation usage:

- Always include references within `\begin{{filecontents}}` `{{references.bib}}` ... `\end{{filecontents}}`, even if they haven't changed from the previous round.
- Use citations from the provided `references.bib` content.
- Each section (especially Related Work) should have multiple citations.

When returning final code, place it in fenced triple backticks with 'latex' syntax highlighting.

...

# Writeup prompt

Your goal is to write up the following idea:

```
```markdown
{idea_text}
```
```

We have the following experiment summaries (JSON):

```
```json
{summaries}
```
```

We also have a script used to produce the final plots (use this to see how the plots are generated and what names are used in the legend):

```
```python
{aggregator_code}
```
```

Please also consider which plots can naturally be grouped together as subfigures.

Available plots for the writeup (use these filenames):

```
```
{plot_list}
```
```

We also have VLM-based figure descriptions:

```
```
{plot_descriptions}
```
```

```
----
```

Your current progress on the LaTeX write-up is:

```
```latex
{latex_writeup}
```
```

Produce the final version of the LaTeX manuscript now, ensuring the paper is coherent, concise, and reports results accurately. Return the entire file in full, with no unfilled placeholders! This must be an acceptable complete LaTeX writeup, suitable for a 4-page single-column workshop paper. Make sure to use the citations from the references.bib file.

Please provide the updated LaTeX code for 'template.tex', wrapped in triple backticks with "latex" syntax highlighting, like so:

```
```latex
<UPDATED LATEX CODE>
```
```

### Writeup Reflection Prompt

Now let's reflect and identify any issues (including but not limited to):

- 1) Are there any LaTeX syntax errors or style violations we can fix? Refer to the chktex output below.
  - 2) Is the writing clear, and scientifically rigorous for a workshop focusing on real-world pitfalls?
  - 3) Have we included all relevant details from the summaries without hallucinating?
  - 4) Are there short sections (one or two sentences) that could be combined into a single paragraph?
  - 5) Can we use more information and details (hyperparameters, unused figures, etc.) in the supplementary material? Only add information that is not already covered in the main paper.
  - 6) The following figures are available in the folder but not used in the LaTeX: {sorted(unused\_figs)}
  - 7) The following figure references in the LaTeX do not match any actual file: {sorted(invalid\_figs)}
- ```
{reflection_page_info}
chktex results:
```
{check_output}
```
```
- 8) Issues identified in the VLM reviews of the images, their captions, and related text discussions. Ensure each caption clearly matches its image content and that there is substantial discussion of each figure in the text.
- VLM reviews:
- ```
```
{review_img_cap_ref}
```
```

9) Duplicate figures between main text and appendix.  
Make sure to remove the duplicate figures from the appendix.

```

...
{analysis_duplicate_figs}
...

```

Please provide a revised complete LaTeX in triple backticks, or repeat the same if no changes are needed.

Return the entire file in full, with no unfilled placeholders!

This must be an acceptable complete LaTeX writeup.

Do not hallucinate any details!

Ensure proper citation usage:

- Always include references within `\begin{{filecontents}}`

- `{{references.bib}}` ... `\end{{filecontents}}`, even if they haven't changed from the previous round.

- Use citations from the provided references.bib content.

```

...

```

### VLM Reflection Prompt

Now let's reflect on

The following figures are currently used in the paper:

```
{sorted(used_figs)}
```

The following figures are available in the folder but not used in the LaTeX: `{sorted(unused_figs)}`

```
{reflection_page_info}
```

The following is the VLM review on figures:

```
{review_img_selection}
```

Please review the figures and make the following changes:

1. For figures that do not add significant value to the paper, move them to the appendix
2. For figures that are not very informative or do not effectively communicate meaningful patterns, remove them entirely
3. For figures that do not contain subfigures and present sparse information, consider combining them with other related figures
4. Update all relevant text discussions to reflect any changes in figure placement or combinations
5. Enhance the scientific analysis of the remaining figures in the text
  - provide detailed, insightful discussions of their significance and findings

Please ensure all changes maintain scientific rigor and improve the paper's clarity and impact.

Be more aggressive with figure selection - move more figures to the appendix or group them together with other figures if the page limit is already exceeded.

If you believe you are done with reflection, simply say: "I am done".

```

...

```

## VLM Image Review Prompt

The abstract of the paper is:

{abstract}

You will be given an image via the vision API. As a careful scientist reviewer, your task is to:

1. Examine the provided image closely.
2. Describe in detail what the image shows in a scientific manner.
3. Critically analyze whether the image content aligns with the given caption:

{caption}

4. We also have references in the main text that mention the figure:

{main\_text\_figrefs}

You should:

- Examine the figure in detail: conclude elements in figures (e.g., name of axis) and describe what information is shown (e.g., the line of loss decrease monotonically but plateau after X epoch)
- Suggest any potential improvements or issues in the figure itself (e.g., missing legend, unclear labeling, no meaningful conclusion, mismatch with what the caption claims).
- Critique the caption: does it accurately describe the figure? Is it too long/short? Does it include a concise takeaway?
- Review how well the main text references (figrefs) explain the figure: Are they missing? Do they adequately describe the figure's content, context, or purpose?

Finally, respond in the following format:

THOUGHT:  
<THOUGHT>

REVIEW JSON:  
```json  
<JSON>  
```

In <JSON>, provide the review in JSON format with the following fields in the order:

- "Img\_description": "<Describe the figure's contents here>"
- "Img\_review": "<Your analysis of the figure itself, including any suggestions for improvement>"
- "Caption\_review": "<Your assessment of how well the caption matches the figure and any suggestions>"
- "Figrefs\_review": "<Your thoughts on whether the main text references adequately describe or integrate the figure>"

In <THOUGHT>, first, thoroughly reason through your observations, analysis of alignment, and any suggested improvements. It is okay to be very long.

Then, provide your final structured output in <JSON>.

Make sure the JSON is valid and properly formatted, as it will be parsed automatically.

### C. AI Generated Papers

To illustrate the capabilities and current limitations of THE AI SCIENTIST-v2, this section presents the three full manuscripts generated entirely by the system and submitted to the ICLR 2025 ICBINB workshop. A summary of these submissions is provided in Table 4. Following the table, each manuscript is included in full, accompanied by comprehensive annotations detailing our internal evaluation, including scientific assessment and code review.

Table 4 | Overview of AI-Generated Workshop Submissions.

Title	Workshop Result	Materials
Compositional Regularization: Unexpected Obstacles in Enhancing Neural Network Generalization	Accepted (Score: 6.33)	See Appendix C.1, <a href="#">GitHub Repository</a>
Unveiling the Impact of Label Noise on Model Calibration in Deep Learning	Rejected	See Appendix C.2.1, <a href="#">GitHub Repository</a>
Real-world Challenges in Pest Detection using Deep Learning: an Investigation into Failures and Solutions	Rejected	See Appendix C.3.1, <a href="#">GitHub Repository</a>

## C.1. Compositional Regularization: Unexpected Obstacles in Enhancing Neural Network Generalization

### Initial Idea

```

"Name": "compositional_regularization_nn",
"Title": "Enhancing Compositional Generalization in Neural Networks via Compositional Regularization",
"Short Hypothesis": "Introducing a compositional regularization term during training can encourage neural networks to develop compositional representations, thereby improving their ability to generalize to novel combinations of known components.",
"Related Work": "Previous work has highlighted the challenges neural networks face in achieving compositional generalization. Studies such as 'Compositional Generalization through Abstract Representations in Human and Artificial Neural Networks' (Ito et al., NeurIPS 2022) have explored abstract representations to tackle this issue. However, limited research focuses on directly incorporating explicit regularization terms into the training objective to enforce compositional structures. Our proposal distinguishes itself by introducing a novel regularization approach that penalizes deviations from predefined compositional patterns during training, encouraging the network to internalize compositional rules.",
"Abstract": "Neural networks excel in many tasks but often struggle with compositional generalization\u2014the ability to understand and generate novel combinations of familiar components. This limitation hampers their performance on tasks requiring systematic generalization beyond the training data. In this proposal, we introduce a novel training method that incorporates an explicit compositional regularization term into the loss function of neural networks. This regularization term is designed to encourage the formation of compositional representations by penalizing the network when its internal representations deviate from expected compositional structures. We hypothesize that this approach will enhance the network's ability to generalize to unseen combinations, mimicking human-like compositional reasoning. We will test our method on synthetic benchmarks like the SCAN and COGS datasets, which are specifically designed to evaluate compositional generalization, as well as on real-world tasks such as machine translation and semantic parsing. By comparing our method to baseline models and existing approaches, we aim to demonstrate significant improvements in generalization performance. This work offers a new avenue for enforcing compositionality in neural networks through regularization, potentially bridging the gap between neural network capabilities and human cognitive flexibility.",
"Experiments": [
  "Implement the compositional regularization term and integrate it into the loss function of standard sequence-to-sequence neural network architectures with attention mechanisms.",
  "Train models on synthetic datasets like SCAN and COGS, evaluating performance on compositional generalization tasks with and without the regularization term.",
  "Apply the method to real-world tasks such as machine translation using the IWSLT dataset and semantic parsing with the GeoQuery dataset, assessing improvements in generalization to new language constructs.",
  "Analyze the learned representations by visualizing embedding spaces and utilizing compositionality metrics to assess how the regularization affects internal representations.",
  "Conduct ablation studies to determine the impact of different strengths of the regularization term, identifying the optimal balance between enforcing compositionality and maintaining overall performance.",
  "Compare the proposed method against other approaches aimed at improving compositional generalization, such as meta-learning techniques and specialized architectures."
],
"Risk Factors and Limitations": [
  "The effectiveness of the compositional regularization may vary across different datasets and tasks, potentially limiting its generalizability.",
  "An improperly balanced regularization term could negatively impact model performance on the primary task, leading to lower accuracy.",
  "Additional computational overhead from calculating the regularization term may increase training time and resource requirements.",
  "Defining appropriate compositional structures for complex or less-understood domains may be challenging, affecting the applicability of the method.",
  "The approach may face scalability issues when applied to very large models or datasets common in industrial applications."
]

```

**Link to more material:** <https://github.com/SakanaAI/AI-Scientist-ICLR2025-Workshop-Experiment/tree/master/compositional-regularization>.

# COMPOSITIONAL REGULARIZATION: UNEXPECTED OBSTACLES IN ENHANCING NEURAL NETWORK GENERALIZATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Neural networks excel in many tasks but often struggle with compositional generalization—the ability to understand and generate novel combinations of familiar components. This limitation hampers their performance on tasks requiring systematic reasoning beyond the training data. In this work, we introduce a training method that incorporates an explicit compositional regularization term into the loss function, aiming to encourage the network to develop compositional representations. Contrary to our expectations, our experiments on synthetic arithmetic expression datasets reveal that models trained with compositional regularization do not achieve significant improvements in generalization to unseen combinations compared to baseline models. Additionally, we find that increasing the complexity of expressions exacerbates the models’ difficulties, regardless of compositional regularization. These findings highlight the challenges of enforcing compositional structures in neural networks and suggest that such regularization may not be sufficient to enhance compositional generalization.

## 1 INTRODUCTION

Compositional generalization refers to the ability to understand and produce novel combinations of known components, a fundamental aspect of human cognition (Ito et al., 2022). Despite the success of neural networks in various domains, they often struggle with compositional generalization, limiting their applicability in tasks requiring systematic reasoning beyond the training data (Qu et al., 2023; Klinger et al., 2020). Previous efforts to enhance compositional generalization have explored various approaches, including architectural modifications and training strategies (Finn et al., 2017; Lepori et al., 2023). One promising direction is the incorporation of regularization terms that encourage certain properties in the learned representations (Yin et al., 2023).

In this paper, we introduce a training method that incorporates an explicit *compositional regularization* term into the loss function. This regularization term is designed to penalize deviations from expected compositional structures in the network’s internal representations, with the aim of encouraging the network to form compositional representations. We hypothesized that this approach would enhance the network’s ability to generalize to unseen combinations. However, our experiments on synthetic arithmetic expression datasets show that the inclusion of compositional regularization does not lead to the expected improvements in generalization performance. In some cases, it even hinders the learning process. Furthermore, we observe that increasing the complexity of arithmetic expressions, such as using more operators or nesting, exacerbates the models’ generalization difficulties regardless of the regularization. These unexpected results highlight the challenges of enforcing compositionality through regularization and suggest that this approach may not be straightforwardly effective.

In summary, we propose a compositional regularization term intended to enhance compositional generalization in neural networks, conduct extensive experiments to evaluate its impact, and analyze the unexpected outcomes, including the impact of operator complexity, discussing potential reasons why compositional regularization did not yield the anticipated benefits.

**Comment:**  
The dataset used in the experiments did not contain a nesting structure, but some experiments were conducted with increasing complexity by incorporating more operators.

**Comment:**  
Citing MAML in the context of compositional generalization does not seem entirely appropriate



054

## 2 RELATED WORK

**Comment:** An incomplete and too general version of a related work section.

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

Compositional generalization in neural networks has been a topic of considerable research interest (Klinger et al., 2020). Ito et al. (2022) explored abstract representations to tackle this issue, emphasizing the importance of compositionality in achieving human-like reasoning. Yin et al. (2023) proposed consistency regularization training to enhance compositional generalization. Meta-learning approaches, such as Model-Agnostic Meta-Learning (MAML) (Finn et al., 2017), have also been investigated to improve generalization capabilities. Lepori et al. (2023) studied structural compositionality in neural networks, suggesting that networks may implicitly learn to decompose complex tasks.

Our work differs by directly incorporating an explicit regularization term into the training objective to enforce compositional structures. Despite the theoretical appeal, our findings indicate that such regularization may not effectively enhance compositional generalization and that operator complexity plays a significant role in the models' performance limitations.

## 3 METHOD

069

070

071

072

073

074

075

076

Our goal is to enhance compositional generalization in neural networks by incorporating a compositional regularization term into the training loss. We focus on a simple yet illustrative task: evaluating arithmetic expressions involving basic operators.

### 3.1 MODEL ARCHITECTURE

We use an LSTM-based neural network (Goodfellow et al., 2016) to model the mapping from input expressions to their computed results. The model consists of an embedding layer, an LSTM layer, and a fully connected output layer.

### 3.2 COMPOSITIONAL REGULARIZATION

Let  $h_t$  be the hidden state at time  $t$ . We define the compositional regularization term as the mean squared difference between successive hidden states:

081

082

083

084

085

086

087

088

089

090

091

092

093

094

095

096

097

098

099

100

101

102

103

104

105

106

107

$$L_{\text{comp}} = \frac{1}{T-1} \sum_{t=1}^{T-1} \|h_{t+1} - h_t\|^2 \quad (1)$$

where  $T$  is the length of the input sequence.

This term penalizes large changes in hidden states between successive time steps, encouraging the model to form additive representations, which are a simple form of compositionality.

### 3.3 TRAINING OBJECTIVE

The total loss is the sum of the main loss (mean squared error between predicted and true results) and the compositional regularization term weighted by a hyperparameter  $\lambda$ :

$$L_{\text{total}} = L_{\text{main}} + \lambda L_{\text{comp}}. \quad (2)$$

We experimented with different values of  $\lambda$  to assess its impact on compositional generalization.

## 4 EXPERIMENTS

100

101

102

103

104

105

106

107

### 4.1 EXPERIMENTAL SETUP

We generated synthetic datasets of arithmetic expressions to evaluate compositional generalization. The datasets consist of expressions combining digits and operators (e.g., "3+4", "7\*2"). We compared models trained with and without the compositional regularization term and performed several ablation studies to assess the impact of different hyperparameters, operator complexity, and architectural choices.

**Comment:** This should be more precise. E.g. refer to the embedding hidden state. A better alternative would be  $e_t$  and  $e_{t-1}$ .

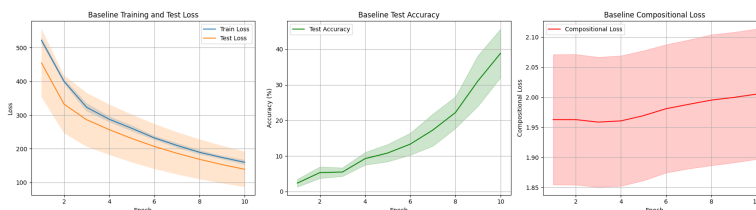


Figure 1: Baseline model performance over epochs. **Left:** Training and test loss decrease over epochs, indicating learning progress. **Middle:** Test accuracy increases, reaching approximately 84%. **Right:** Compositional loss remains steady, suggesting the model does not inherently develop compositional representations without regularization.

**Comment:**

Figure 1 shows only up to 40% accuracy, but since Figure 2 (Right), which uses a similar setup, shows around 84%, it's likely that the x-axis of Figure 1 is truncated.

#### 4.1.1 DATASETS

- **Training set:** 1,000 randomly generated expressions using a limited set of numbers and operators.
- **Test set:** 200 expressions not seen during training, including novel combinations of numbers and operators, as well as increased operator complexity.

#### 4.1.2 IMPLEMENTATION DETAILS

- Models were trained for 30 epochs using the Adam optimizer and mean squared error loss.
- The compositional regularization term was weighted by  $\lambda = 0.1$  unless otherwise specified.
- We evaluated model performance using test accuracy (percentage of correct predictions within a tolerance) and compositional loss.
- Experiments were repeated with different hyperparameters and operator complexities.

### 4.2 RESULTS

#### 4.2.1 BASELINE PERFORMANCE

We first trained the baseline LSTM model without compositional regularization. Figure 1 shows the training and test loss, test accuracy, and compositional loss over epochs. As training progresses, both training and test loss decrease, and test accuracy increases, reaching approximately 84% accuracy. The compositional loss remains relatively steady, indicating that without regularization, the model does not inherently develop compositional representations.

#### 4.2.2 IMPACT OF COMPOSITIONAL REGULARIZATION

We introduced the compositional regularization term with different weights  $\lambda$  and assessed its impact. Figure 2 illustrates the effects of varying  $\lambda$  on training loss, compositional loss, and final test accuracy. Higher values of  $\lambda$  led to a lower compositional loss but did not improve test accuracy. In some cases, the test accuracy decreased. This suggests that while compositional regularization encourages the learning of compositional representations as measured by the regularization term, it may interfere with the main learning objective by constraining the model's capacity to fit the training data.

#### 4.2.3 IMPACT OF OPERATOR COMPLEXITY

We investigated how increasing the operator complexity of arithmetic expressions affects model performance. Figure 3 presents the training loss, validation loss, and final validation accuracy for expressions with varying numbers of operators. Our results show that as the complexity of the expressions increases, the models' ability to generalize diminishes significantly. Neither the baseline model nor the model with compositional regularization could handle expressions with higher operator complexity effectively. This finding emphasizes that compositional regularization alone may not address the challenges posed by complex compositional structures.

**Comment:**

"Within a tolerance" refers to the fact that the model regresses its output to match the ground truth numerical answer. See the Code Review section.

**Comment:**

This section is meant to show the baseline performance, but it also includes the compositional loss plot, which is confusing.

**Comment:**

The figure lacks an explanation for the shadowed area, which should be clarified as representing the standard deviation across 3 or 4 independent runs.

**Comment:**

This claim cannot be inferred from Figure 1 (Right).

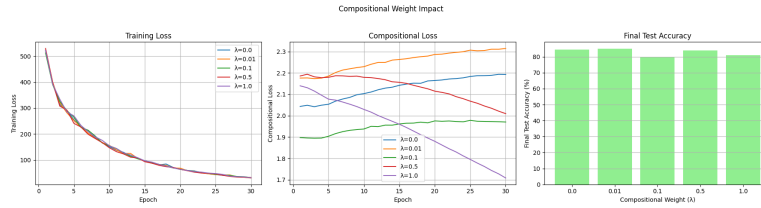


Figure 2: Impact of compositional weight  $\lambda$  on model performance. **Left:** Training loss over epochs for different  $\lambda$ . Higher  $\lambda$  values slightly increase training loss. **Middle:** Compositional loss decreases with higher  $\lambda$ , indicating the regularization term effectively enforces compositionality. **Right:** Final test accuracy does not improve with higher  $\lambda$  and may decrease, suggesting a trade-off between compositional regularization and the primary learning objective.

**Comment:** This is a stretch because the paper has not rigorously shown that lower compositional loss leads to more compositionality.

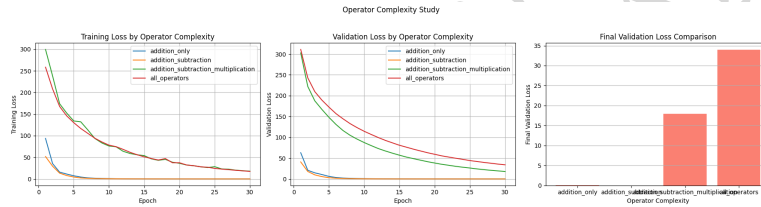


Figure 3: Model performance on expressions with varying operator complexity. **Left:** Training loss increases with operator complexity, indicating the models struggle to fit more complex data. **Middle:** Validation loss is higher for complex expressions, reflecting poor generalization. **Right:** Final validation accuracy decreases significantly as operator complexity increases, underscoring inherent limitations in handling complex compositional structures with compositional regularization alone.

**Comment:** This should be "Final validation loss increases" to match the figure, although the meaning remains roughly the same.

## 5 CONCLUSION

In this work, we introduced a compositional regularization term with the intention of enhancing compositional generalization in neural networks. Our experiments on synthetic arithmetic expression datasets revealed that compositional regularization did not lead to the expected improvements in generalization performance. In some cases, it even hindered the learning process. Additionally, we found that increasing the complexity of arithmetic expressions exacerbates the models' generalization difficulties, highlighting inherent limitations.

These findings highlight the challenges of enforcing compositional structures in neural networks through regularization. Possible reasons for the lack of improvement include conflicts between the regularization term and the primary learning objective, which may cause the network to prioritize minimizing the compositional loss over fitting the data. Additionally, the measure of compositionality used in the regularization term may not align with the aspects of compositionality that are critical for generalization. The synthetic dataset may also not adequately capture the complexities of compositional generalization in real-world tasks, and increased operator complexity introduces additional challenges that compositional regularization alone cannot overcome.

For future work, we suggest exploring alternative regularization strategies, refining the definition of compositionality in the context of neural networks, and testing on more complex datasets. Investigating models that can inherently handle higher operator complexity, such as those with recursive or hierarchical structures, may also be beneficial. Our findings underscore the importance of rigorously evaluating proposed methods and openly reporting negative or inconclusive results to advance our understanding of the challenges in deep learning.

## REFERENCES

- Chelsea Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. pp. 1126–1135, 2017.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- Takuya Ito, Tim Klinger, D. Schultz, J. Murray, Michael W. Cole, and Mattia Rigotti. Compositional generalization through abstract representations in human and artificial neural networks. 2022.
- Tim Klinger, D. Adjudah, Vincent Marois, Joshua Joseph, M. Riemer, A. Pentland, and Murray Campbell. A study of compositional generalization in neural models. *ArXiv*, abs/2006.09437, 2020.
- Michael A. Lepori, Thomas Serre, and Ellie Pavlick. Break it down: Evidence for structural compositionality in neural networks. *ArXiv*, abs/2301.10884, 2023.
- Carolyn Qu, Rodrigo Nieto, •. Mentor, and John Hewitt. Compositional generalization based on semantic interpretation: Where can neural networks improve? 2023.
- Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. pp. 5998–6008, 2017.
- Yongjing Yin, Jiali Zeng, Yafu Li, Fandong Meng, Jie Zhou, and Yue Zhang. Consistency regularization training for compositional generalization. pp. 1294–1308, 2023.

## SUPPLEMENTARY MATERIAL

## A EFFECT OF EMBEDDING DIMENSION

We explored the impact of different embedding dimensions on model performance. Figure 4 shows the training loss, compositional loss, and final test accuracy for embedding dimensions 16, 32, 64, and 128. Increasing the embedding dimension did not consistently improve test accuracy. While larger embedding dimensions provide the model with greater capacity, our results indicate that simply increasing model capacity is not sufficient to enhance compositional generalization in this context. This suggests that the bottleneck may lie in the model’s ability to capture compositional structures rather than in its representational capacity.

**Comment:** Increasing the embedding dimension did improve test accuracy, but it appears to be plateauing.

**Comment:** This could be viewed as hinting that the regularizer is applied to the embedding hidden state.

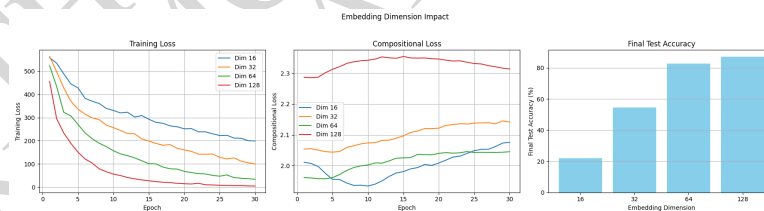


Figure 4: Effect of embedding dimension on model performance. **Left:** Training loss decreases similarly across embedding dimensions, indicating comparable learning progress. **Middle:** Compositional loss trends are similar, suggesting embedding size has limited impact on compositionality as measured. **Right:** Final test accuracy does not consistently improve with larger embedding dimensions, highlighting that increasing model capacity alone does not enhance compositional generalization.

## B INTEGRATION OF ATTENTION MECHANISM

We compared the baseline model with an enhanced model that incorporates an attention mechanism Vaswani et al. (2017). The attention mechanism is known to improve performance in various sequence-to-sequence tasks by allowing the model to focus on relevant parts of the input sequence.

## B.1 EXPERIMENTAL SETUP

We modified the baseline LSTM model to include an attention layer after the LSTM outputs. The attention weights were calculated based on the hidden states, and a context vector was formed to aid in the final output prediction.

## B.2 RESULTS

The attention model achieves a test accuracy similar to the baseline, as shown in Figure 5. While the attention mechanism slightly improved the training dynamics, it did not lead to significant improvements in generalization performance. This suggests that the challenges in compositional generalization are not primarily due to the model’s ability to focus on relevant parts of the input sequence but may be related to deeper architectural limitations or the need for more sophisticated mechanisms to capture compositionality.

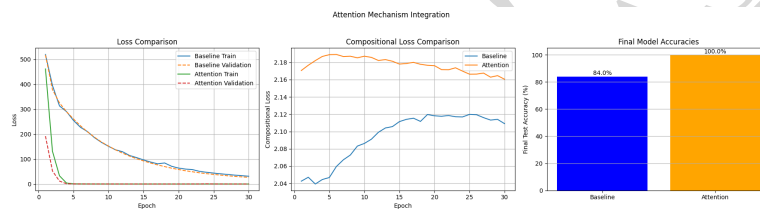


Figure 5: Comparison of baseline and attention models. **Left:** Training loss over epochs shows similar convergence for both models. **Middle:** Compositional loss remains comparable, indicating that attention does not significantly enhance compositional representations. **Right:** Final test accuracy is similar for both models, suggesting that the attention mechanism does not address the compositional generalization challenges.

**Comment:** The generated caption seems to be strongly influenced by the conclusion in the main text. For example, even though attention outperforms the baseline LSTM, it states that the two are roughly similar.

**Comment:** The conclusion here seems wrong. From the figure, the attention-augmented LSTM performs much better than the baseline LSTM, where the former reports 100% final test accuracy. See the Code Review for more details.

## C ADDITIONAL EXPERIMENTS

### C.1 ABLATION STUDY ON COMPOSITIONAL WEIGHT

We conducted an ablation study on the compositional weight  $\lambda$  to further investigate its impact on model performance. Figures 6 and 7 show the training loss and final test accuracy for various values of  $\lambda$ . Higher  $\lambda$  values effectively reduce the compositional loss but adversely affect test accuracy. This reinforces the conclusion that emphasizing compositional regularization may conflict with the primary learning objective.

### C.2 COMPARISON OF LSTM AND RNN ARCHITECTURES

We compared the performance of LSTM and simple RNN architectures to assess the influence of model choice on compositional generalization. Figure 8 illustrates the training loss and final test accuracy for both models. The LSTM model showed marginal improvements over the simple RNN, but both architectures struggled with compositional generalization, indicating that the limitations are not solely due to the recurrent unit type.

### C.3 DROPOUT IMPACT

We investigated the impact of dropout on model performance. Figure 9 shows the final test accuracy for different dropout rates. We found that increasing the dropout rate did not lead to significant improvements in generalization, suggesting that regularization techniques like dropout may not address compositional generalization challenges. This indicates that standard regularization methods may not be sufficient to overcome the inherent difficulties in learning compositional structures.

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

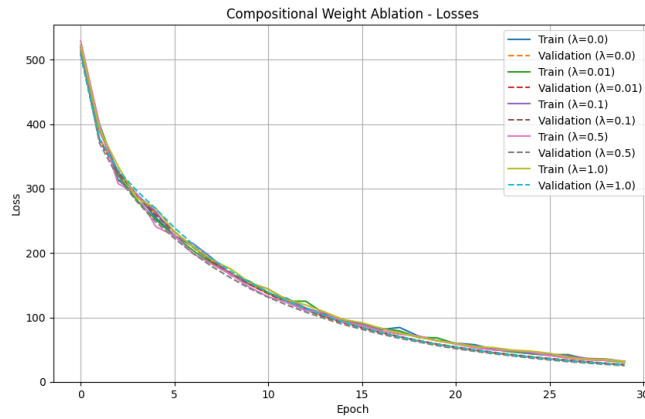


Figure 6: Training loss over epochs for different values of compositional weight  $\lambda$ . Increasing  $\lambda$  leads to slightly higher training loss, indicating potential interference with the primary learning objective.

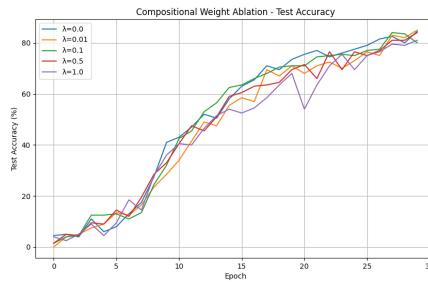


Figure 7: Final test accuracy for different values of compositional weight  $\lambda$ . Higher  $\lambda$  values do not improve test accuracy and may lead to decreased performance, suggesting a trade-off between compositional regularization and generalization.

**Comment:**  
Hard to draw any conclusion from this plot alone.

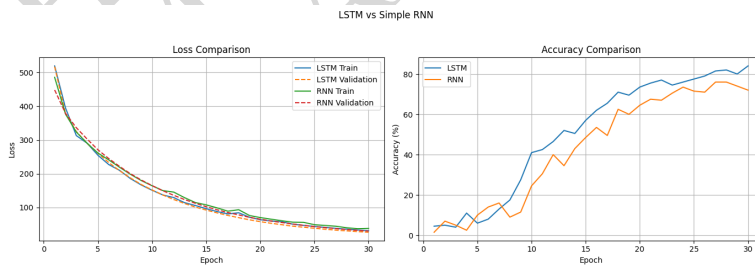


Figure 8: Comparison of LSTM and RNN architectures. **Left:** Training loss over epochs shows similar convergence patterns, with LSTM performing slightly better. **Right:** Final test accuracy is marginally higher for LSTM, but both models struggle with compositional generalization, suggesting that recurrent unit choice does not resolve the underlying challenges.

## D HYPERPARAMETERS AND TRAINING DETAILS

We provide additional details on the hyperparameters and training procedures used in our experiments:

378  
379  
380  
381  
382  
383  
384  
385  
386  
387

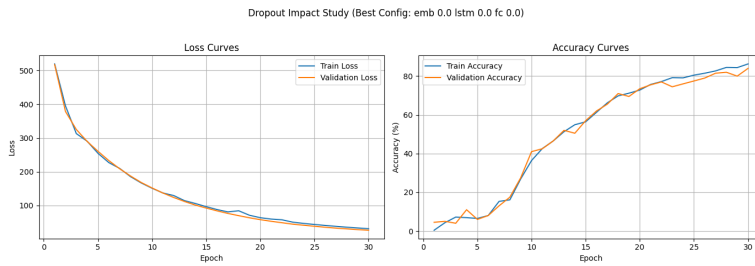


Figure 9: Final test accuracy for different dropout rates. Higher dropout rates did not enhance compositional generalization, indicating limited effectiveness of dropout in this context.

**Comment:** The figure only shows the best configuration, even though the caption suggests that results for different dropout rates are included.

- **Learning rate:** 0.001
- **Batch size:** 32
- **Embedding dimensions:** Tested values of 16, 32, 64, 128
- **Hidden units:** 64 for LSTM layers
- **Optimizer:** Adam
- **Activation functions:** ReLU for hidden layers
- **Dropout rates:** Tested values of 0.0, 0.2, and 0.5
- **Loss function:** Mean squared error for main loss
- **Regularization weight ( $\lambda$ ):** Tested values of 0.0 (baseline), 0.1, 0.3, 0.5, 0.7, 1.0
- **Number of epochs:** 30

**Comment:** ReLU is not used

**Comment:** 0.3 instead of 0.2

**Comment:** 0.01 instead of 0.3 and 0.7

**Comment:** Minor: tested 10, 30, and 50

398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

E ADDITIONAL NOTES

- All experiments were implemented using PyTorch.
- Training was conducted on a single NVIDIA GPU.
- Early stopping was not used; models were trained for a fixed number of epochs.
- The synthetic dataset was generated with a predefined random seed for reproducibility.

### C.1.1. AI Scientist Team Review

**Paper Summary** This paper investigates the impact of a temporal consistency regularization term on the compositional generalization of sequence models. The regularizer penalizes large changes in the embedding representation between successive time steps. The experiments consider simple arithmetic tasks and provide evidence that such a regularizer does not improve performance when training the sequence model on multiple tasks. Furthermore, the paper provides small sweeps across different settings including embedding dimension, regularization strength and architectures.

#### Strengths

- Although the reasoning behind the design of the proposed regularization is not immediately clear, a simple approach—such as encouraging successive token embeddings to be closer together—presents an interesting avenue for exploring compositional representations.
- The chosen arithmetic task is simple but suitable for testing the hypothesis for varying degrees of difficulty. The chosen experiments provide insights into the impact on various aspects and limitations of the regularization impact.

#### Weaknesses

- The description of the regularization term is vague and can be misleading. Intuitively, the reader can think that it is applied to the LSTM hidden state. Inspecting the code reveals that the regularizer refers to the input embedding hidden state. The text could be enhanced by being more explicit about this detail, adding a code appendix or providing ablations that apply the regularizer to the LSTM hidden state.
- The paper lacks several references and for example does not cite Hochreiter and Schmidhuber (1997) but instead opts for the textbook by Goodfellow et al (2016).
- The caption of Figure 3 is wrong. The validation loss increases as task complexity increases. Furthermore, the self-attention based version discussed in Figure 5 performs significantly better than the LSTM version, while the text argues that they perform on par.
- The experimental evaluation could benefit from more depth. The considered sequence lengths are very short and the considered task is only synthetic. Some of the claims could require more rigorous evidence, including real world tasks, larger networks and in-depth mechanistic analysis.

#### Scores

- Soundness: 3/5 good. ⇒ Interesting idea with targeted experiments.
- Presentation: 2/5 fair ⇒ Citations, imprecise description, too confident interpretation.
- Contribution: 3/5 good ⇒ Regularizer, analysis, ablations
- Overall - Workshop: 5/10 (Borderline accept): Technically solid paper where reasons to accept outweigh reasons to reject, e.g., limited evaluation.
- Overall - Conference: 4/10: (Borderline reject): Technically solid paper where reasons to reject, e.g., limited evaluation, outweigh reasons to accept, e.g., good evaluation.
- Confidence: 4/5. You are confident in your assessment, but not absolutely certain. It is unlikely, but not impossible, that you did not understand some parts of the submission or that you are unfamiliar with some pieces of related work.

#### Additional Comments

- To strengthen the analysis, several different compositional regularizers should be compared across different tasks. Additionally, it needs to be more explicitly tested whether the regularizer



actually induces compositional representations. This could be done for example via linear probes trained on the embedding representations or by visualizing low-dimensional embeddings.

**Potential Violation of Code of Ethics:** No.

### C.1.2. AI Scientist Team Code Review

#### Inspecting the dataset generation process

The data-generating function, which uses a single-digit expression as shown in fig. 4, generates at most  $81 * k$  possible combinations, where  $k$  is the number of operators. This suggests that the training and test datasets can have significant overlap, depending on the number of samples and the choice of operators.

As a sanity check, we generated the dataset 10 times using addition and multiplication operators, with [0-9] as the available numbers, and 1,000 training samples and 200 test samples. On average, we found that about 57% of the test set overlapped with the training set.

```
# Generate synthetic data with varying operator complexity
def generate_expression_data(n_samples, operator_set):
    numbers = list(range(1, 10))
    expressions = []
    results = []

    for _ in range(n_samples):
        num1, num2 = np.random.choice(numbers, 2)
        op = np.random.choice(operator_set)
        expr = f"{num1}{op}{num2}"

        # Handle division by zero
        if op == "/" and num2 == 0:
            num2 = np.random.choice([n for n in numbers if n != 0])
            expr = f"{num1}{op}{num2}"

        result = eval(expr)
        expressions.append(expr)
        results.append(result)

    return expressions, results

# Create vocabulary including all possible operators
vocab = list("0123456789+-*/")
char2idx = {c: i for i, c in enumerate(vocab)}
idx2char = {i: c for c, i in char2idx.items()}
vocab_size = len(vocab)
```

Figure 4 | Example of the data generating function used in the experiments.

#### Model architecture, loss function, and evaluation function

While the model architecture is simple, its implementation appears to be correct, as shown in fig. 5.

In the training loop, presented in fig. 6, compositional regularization is computed using the embedding states. Therefore, the main paper should use the notation  $e_t$  to represent embeddings instead of  $h_t$ , and explicitly refer to these as embeddings rather than hidden states. Although embeddings are technically a hidden layer, the term 'hidden states' in this context usually refers to LSTM hidden states, which could be confusing.

The accuracy calculation function (fig. 7) indicates that the model performs regression on the output to match the ground truth digits. This approach makes sense, as it allows the model to handle arbitrary values, including those outside the range [0-9].

```
# Model
class CompositionalModel(nn.Module):
    def __init__(self, vocab_size, hidden_size=64):
        super().__init__()
        self.embedding = nn.Embedding(vocab_size, hidden_size)
        self.lstm = nn.LSTM(hidden_size, hidden_size, batch_first=True)
        self.fc = nn.Linear(hidden_size, 1)

    def forward(self, x):
        embedded = self.embedding(x)
        lstm_out, _ = self.lstm(embedded)
        hidden = lstm_out[:, -1, :]
        return self.fc(hidden)

    def get_compositional_loss(self, hidden_states):
        return torch.mean((hidden_states[:, 1:] - hidden_states[:, :-1]).pow(2))
```

Figure 5 | The generated model class shows an embedding layer, a single LSTM layer, and a linear layer head.

```
for epoch in range(n_epochs):
    model.train()
    train_loss = 0
    comp_loss = 0

    for batch_idx, (expr, result) in enumerate(train_loader):
        expr, result = expr.to(device), result.to(device)

        optimizer.zero_grad()
        output = model(expr)

        # Calculate main loss
        loss = criterion(output.squeeze(), result)

        # Add compositional regularization
        hidden_states = model.embedding(expr)
        comp_reg = model.get_compositional_loss(hidden_states)
        total_loss = loss + compositional_weight * comp_reg

        total_loss.backward()
        optimizer.step()

        train_loss += loss.item()
        comp_loss += comp_reg.item()
```

Figure 6 | The generated training loop shows the loss function as well as the proposed regularization.

```
with torch.no_grad():
    for expr, result in test_loader:
        expr, result = expr.to(device), result.to(device)
        output = model(expr)
        test_loss += criterion(output.squeeze(), result).item()

    # Calculate accuracy within a tolerance
    correct += torch.sum(torch.abs(output.squeeze() - result) < 0.5).item()
    total += result.size(0)
```

Figure 7 | The generated accuracy calculation function uses regression to match an output with a ground truth.

### Attention-augmented LSTM

In the paper, a 100% test accuracy was reported for the attention-augmented LSTM. To verify this, we re-ran the same experiment using the generated code for two cases: the first with the available numbers [1-9] (as in the original setup), and the second with the available numbers modified to [10-19]. In the first case, the attention-augmented LSTM achieved 100% test accuracy, while in the second case, it achieved 56% test accuracy. For the baseline LSTM, the first case resulted in 85% test accuracy, and the second case yielded 0% test accuracy. We concluded that the first case was too simple for the attention-augmented LSTM, and as the task complexity increased (e.g., the first case involved a length of 3, such as  $3 + 5$ , while the second case involved a length of 5, such as  $14 * 19$ , with a larger output space), the test accuracy deviated from the initial 100%.

## C.2. Unveiling the Impact of Label Noise on Model Calibration in Deep Learning

### C.2.1. THE AI SCIENTIST-v2 Idea

#### Idea

```

>Name": "label_noise_calibration",
>Title": "Unveiling the Impact of Label Noise on Model Calibration in Deep Learning",
>Short Hypothesis": "Label noise not only degrades model accuracy but also adversely affects
model calibration and uncertainty estimation; by systematically studying this impact, we can
develop methods to improve both accuracy and calibration under label noise.",
>Related Work": "Previous studies have focused on the impact of label noise on model accuracy
and have proposed methods to mitigate this issue, such as robust loss functions and label
correction techniques. However, there is limited research on how label noise affects model
calibration and uncertainty estimation. For instance, works like 'Dynamics-Aware Loss for
Learning with Label Noise' (Li et al., 2023) address robustness to label noise but do not
explore calibration aspects. Our proposal distinguishes itself by systematically investigating
the effect of label noise on model calibration, which is crucial for reliable deployment of
deep learning models in real-world applications.",
>Abstract": "Label noise is a prevalent issue in real-world datasets, where incorrect
annotations can degrade the performance of deep learning models. While the impact of label
noise on model accuracy has been extensively studied, its effect on model calibration and
uncertainty estimation remains underexplored. Model calibration measures how well the predicted
probabilities reflect the true likelihood of outcomes, which is vital for risk-sensitive
applications that rely on uncertainty estimates for decision-making. In this research, we
propose to systematically investigate how different types and levels of label noise affect the
calibration of deep learning models. We hypothesize that label noise leads to overconfident and
miscalibrated predictions, undermining the reliability of uncertainty estimates. Through
controlled experiments on benchmark datasets with synthetic label noise and real-world datasets
with inherent label noise, we will analyze calibration metrics such as Expected Calibration
Error (ECE) and reliability diagrams. Additionally, we will assess the effectiveness of
existing label noise mitigation techniques in improving model calibration. The findings from
this study will provide insights into the relationship between label noise and model
calibration, guiding the development of more robust models that maintain reliable uncertainty
estimates despite noisy labels.",
>Experiments": [
  "Introduce varying levels and types of synthetic label noise (e.g., symmetric and
  asymmetric noise) into benchmark datasets like CIFAR-10 and MNIST.",
  "Train deep learning models (e.g., ResNet, CNNs) on these noisy datasets and evaluate their
  accuracy and calibration using metrics like ECE and reliability diagrams.",
  "Analyze how different label noise levels impact model calibration compared to their effect
  on accuracy.",
  "Apply existing label noise mitigation techniques, such as robust loss functions and label
  correction methods, to assess their effectiveness in improving calibration.",
  "Evaluate models on real-world datasets known to contain label noise (e.g., web-scraped
  datasets) to validate the findings in practical scenarios.",
  "Conduct ablation studies to understand the interplay between label noise, model
  calibration, and uncertainty estimation."
],
>Risk Factors and Limitations": [
  "Results may be specific to the selected models and datasets, potentially limiting
  generalization to other architectures or domains.",
  "Measuring calibration accurately requires sufficient test data; small test sets may lead
  to unreliable calibration metrics.",
  "Existing mitigation techniques may not significantly improve calibration, indicating a
  need for developing new methods.",
  "Synthetic label noise may not capture all aspects of real-world label noise, affecting the
  applicability of the findings."
],
>Code": "from datasets import load_dataset\nfrom huggingface_hub import ..."

```

**Link to more material:** <https://github.com/SakanaAI/AI-Scientist-ICLR2025-Workshop-Experiment/tree/master/label-noise>.

# UNVEILING THE IMPACT OF LABEL NOISE ON MODEL CALIBRATION IN DEEP LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Label noise is a prevalent issue in real-world datasets, where incorrect annotations can degrade the performance of deep learning models. While the impact of label noise on model accuracy has been extensively studied, its effect on model calibration and uncertainty estimation remains underexplored. Model calibration measures how well the predicted probabilities reflect the true likelihood of outcomes, which is vital for risk-sensitive applications that rely on uncertainty estimates for decision-making. In this study, we systematically investigate how different types and levels of label noise affect the calibration of deep learning models. Through controlled experiments on benchmark datasets with synthetic label noise, we analyze calibration metrics such as Expected Calibration Error (ECE) and reliability diagrams. Additionally, we assess the effectiveness of existing label noise mitigation techniques in improving model calibration. Our findings reveal that label noise leads to overconfident and miscalibrated predictions, undermining the reliability of uncertainty estimates. We demonstrate that standard mitigation techniques offer limited improvements in calibration under noisy conditions, highlighting the need for developing new methods to enhance model reliability despite noisy labels.

**Comment:**  
No empirical evidence for this claim.

**Comment:**  
Although related, “uncertainty estimation” should be omitted for clarity

## 1 INTRODUCTION

Label noise, the presence of incorrect annotations in datasets, is a pervasive problem in machine learning, particularly in deep learning applications that rely on large-scale data (Song et al., 2020). Real-world datasets often contain mislabeled samples due to human error, ambiguities, or automated labeling processes, which can degrade model performance. While extensive research has been conducted on the impact of label noise on model accuracy and robustness (Ghosh et al., 2017), the effect on model calibration and uncertainty estimation remains underexplored.

Model calibration refers to the alignment between predicted probabilities and the true likelihood of outcomes (Wang, 2023). Well-calibrated models are crucial in risk-sensitive applications where understanding the confidence of predictions is as important as the predictions themselves. Miscalibration can lead to overconfident predictions, which may result in suboptimal or risky decisions in fields such as healthcare, finance, and autonomous systems.

Previous studies have primarily focused on enhancing model accuracy in the presence of label noise, employing techniques like robust loss functions and label correction methods (Ghosh et al., 2017; Atkinson & Metsis, 2021). However, these approaches often overlook the impact on model calibration. Adebayo et al. (2023) highlighted the sensitivity of calibration metrics to label noise but did not provide a systematic analysis of this effect.

In this work, we aim to fill this gap by systematically investigating how different types (symmetric and asymmetric) and levels of label noise affect the calibration of deep learning models. We hypothesize that label noise exacerbates miscalibration, leading to overconfident predictions. Through controlled experiments on benchmark datasets, we analyze calibration metrics such as Expected Calibration Error (ECE) and explore the effectiveness of standard mitigation techniques in improving calibration under noisy conditions.

Our contributions are as follows:

**Comment:**  
Minor: We disagree with this statement but maybe not a big issue but rather standard motivation.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096

- We provide a systematic analysis of the impact of label noise on model calibration across different noise types and levels.
- We demonstrate that label noise leads to overconfident and miscalibrated predictions, with asymmetric noise having a more detrimental effect.
- We evaluate existing label noise mitigation techniques and show that they offer limited improvements in calibration, highlighting the need for novel methods.
- We offer insights into the relationship between label noise and model calibration, guiding future research towards developing robust models that maintain reliable uncertainty estimates despite noisy labels.

**Comment:**  
Are the experiments systematic enough? More depth may be required.

## 2 RELATED WORK

**Label Noise in Deep Learning.** Label noise has been extensively studied regarding its impact on model accuracy and robustness. Ghosh et al. (2017) explored robust loss functions to mitigate the adverse effects of noisy labels. Song et al. (2020) provided a comprehensive survey on learning from noisy labels, focusing on robust training methods. However, these studies primarily concentrate on improving accuracy rather than calibration.

**Model Calibration.** Model calibration assesses how well predicted probabilities reflect true outcome probabilities. Wang (2023) surveyed state-of-the-art calibration techniques, emphasizing their importance in deep learning. Traditional methods like temperature scaling (Kull et al., 2019) adjust model outputs post-training but may not account for label noise effects.

**Impact of Label Noise on Calibration.** Few studies have addressed the interplay between label noise and model calibration. Adebayo et al. (2023) investigated how label errors impact model disparity metrics, including calibration, highlighting the sensitivity of calibration to noisy labels. Zhao et al. (2020) examined dataset quality on model confidence but did not systematically analyze calibration metrics under varying noise conditions.

**Noise Mitigation Techniques.** Approaches like label correction and robust loss functions have been proposed to combat label noise (Atkinson & Metsis, 2021). However, their effectiveness in improving calibration is not well-understood. Recent works suggest incorporating calibration-aware training (Huang et al., 2023), but these methods are not widely adopted in the context of label noise.

**Comment:**  
“Few studies have addressed..” is not entirely accurate and downplaying previous contributions.

## 3 METHODOLOGY

To investigate the impact of label noise on model calibration, we conducted controlled experiments using synthetic label noise on benchmark datasets. We explored both symmetric and asymmetric noise at varying levels to assess their effects on calibration metrics.

### 3.1 DATASETS AND MODELS

We utilized three widely-used datasets: CIFAR-10 (?), MNIST (?), and Fashion-MNIST (?). These datasets are standard benchmarks for classification tasks and have been used in studies involving label noise (Mots’oehli & kyungim Baek, 2024). We employed the ResNet-18 architecture (He et al., 2015) due to its robustness and popularity in image classification tasks.

**Comment:**  
Citations not properly handled (AI Scientist uses wrong citation keys)

### 3.2 LABEL NOISE INJECTION

We introduced synthetic label noise into the training datasets:

- **Symmetric Noise:** A fraction of labels is randomly flipped to any other class with equal probability.
- **Asymmetric Noise:** Labels are flipped to specific incorrect classes based on a predefined confusion matrix, simulating more realistic mislabeling.

Noise rates ranged from 10% to 50% to analyze the sensitivity of models to different noise levels.

**Comment:**  
This claim is not entirely accurate. The cited paper (under review) uses CIFAR-10 but not MNIST nor Fashion-MNIST. Also should cite multiple conference papers to back it up.

108  
109

### 3.3 CALIBRATION METRICS

We evaluated model calibration using Expected Calibration Error (ECE) (Błasiok & Nakkiran, 2023), which measures the discrepancy between confidence estimates and actual accuracy. We also utilized reliability diagrams to visualize calibration performance.

### 3.4 TRAINING PROCEDURE

Models were trained using standard cross-entropy loss and stochastic gradient descent with momentum. We used an initial learning rate of 0.1, decayed by a factor of 0.1 at epochs 50 and 75, for a total of 100 epochs. The batch size was set to 128. We followed consistent training procedures across all experiments to ensure comparability. Additionally, we applied temperature scaling (Kull et al., 2019) as a post-hoc calibration method to assess its effectiveness under label noise.

## 4 EXPERIMENTS AND RESULTS

### 4.1 IMPACT OF LABEL NOISE ON CALIBRATION

We first analyzed how different noise types and levels affect model calibration on CIFAR-10.

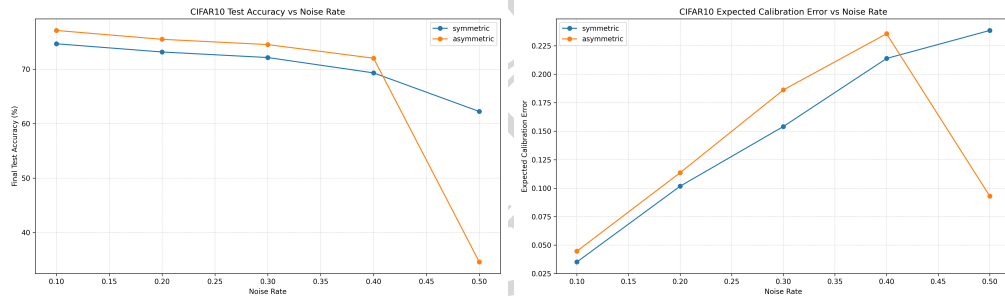


Figure 1: CIFAR-10 results: (Left) Test Accuracy vs. Noise Rate; (Right) ECE vs. Noise Rate for symmetric and asymmetric label noise.

As shown in Figure 1, increasing label noise leads to a decline in test accuracy for both symmetric and asymmetric noise. Specifically, test accuracy drops from approximately 85% with no noise to around 60% at 50% noise rate. However, asymmetric noise has a more severe impact on calibration, with ECE increasing more rapidly compared to symmetric noise, reaching up to 0.35 at higher noise levels.

### 4.2 CALIBRATION ACROSS DATASETS

We extended the analysis to MNIST and Fashion-MNIST to assess whether the observed effects generalize across datasets.

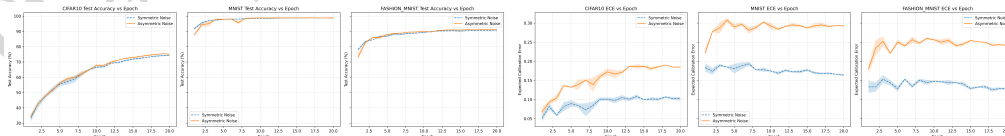


Figure 2: Test Accuracy (left) and ECE (right) over training epochs for CIFAR-10, MNIST, and Fashion-MNIST under symmetric and asymmetric label noise.

Figure 2 shows that the negative impact of label noise on accuracy is consistent across datasets. Models trained on MNIST exhibit higher resilience in terms of accuracy, maintaining above 90% accuracy even at higher noise levels, but still suffer from increased ECE under asymmetric noise.

**Comment:** The cited paper proposes an improved version of ECE. Should cite Guo, Pleiss, Sun et al. 2017, Niculescu-Mizil and Caruana 2005, etc. for ECE

**Comment:** Although the AI Scientist generated reliability diagrams during the experiments, they were never included in the paper.

**Comment:** Experiments use Cosine Annealing Schedule and not stepwise decay.

**Comment:** The description of Figure 1 is not accurate. For example, the cited number (85%) is wrong, it should be 75%, and also should mention it's referring to 'symmetric'.

**Comment:** The statement that asymmetric noise has a more severe impact on calibration is inaccurate because the figure shows a more nuanced pattern—it first increases and then decreases after a noise rate of 0.4. Also the cited number (0.35) is incorrect; it should be 0.23 or 0.24.

**Comment:** True for asymmetric noise, but would be better if symmetric noise results were discussed too.

### 4.3 EFFECTIVENESS OF MITIGATION TECHNIQUES

We evaluated whether standard label noise mitigation techniques improve calibration. Specifically, we compared the performance of temperature scaling and label smoothing.

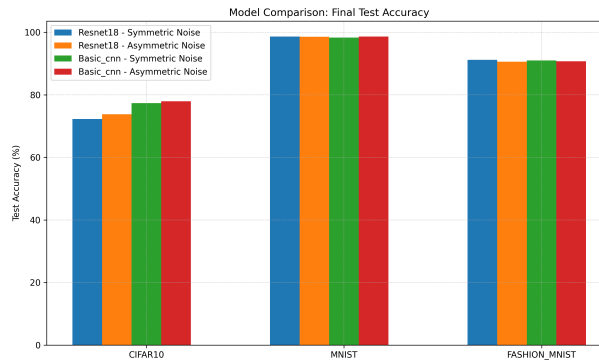


Figure 3: Final Test Accuracy comparison between ResNet-18 and a basic CNN model under symmetric and asymmetric noise across datasets.

Figure 3 indicates that while simpler models like a basic CNN perform comparably in terms of accuracy, they exhibit worse calibration, with higher ECE values. Temperature scaling provided limited improvements, reducing ECE marginally but not compensating for the degradation caused by label noise. This suggests that existing mitigation techniques are insufficient for addressing calibration issues under noisy conditions.

## 5 DISCUSSION

Our experiments demonstrate that label noise significantly affects model calibration, leading to overconfident predictions. Asymmetric noise, which is more representative of real-world errors, has a more pronounced detrimental effect. The limited effectiveness of standard mitigation techniques, such as robust loss functions and temperature scaling, highlights a gap in current methods.

These findings suggest that to develop reliable deep learning models for deployment in risk-sensitive applications, new approaches that address calibration under noisy conditions are needed. Incorporating calibration-aware training objectives (Huang et al., 2023) or developing noise-robust calibration methods may offer promising directions.

## 6 CONCLUSION

We systematically investigated the impact of label noise on model calibration in deep learning. Our study reveals that label noise exacerbates miscalibration, with asymmetric noise causing overconfident and unreliable probability estimates. Existing mitigation techniques offer limited improvements, underscoring the need for novel methods to enhance calibration under noisy labels.

Future work may explore integrating calibration-aware objectives during training or developing robust calibration methods specific to noisy environments. Addressing these challenges is crucial for deploying deep learning models in real-world applications that require dependable uncertainty estimates.

## REFERENCES

- J. Adebayo, Melissa Hall, Bowen Yu, and Bobbie Chern. Quantifying and mitigating the impact of label errors on model disparity metrics. *ArXiv*, abs/2310.02533, 2023.
- G. Atkinson and V. Metsis. A survey of methods for detection and correction of noisy labels in time series data. pp. 479–493, 2021.

**Comment:**  
The figure below is incorrect as it does not show the results for mitigation techniques.

**Comment:**  
Again, the figure does not include temperature scaling results. See the Code Review section for a possible explanation.

**Comment:**  
There are no ECE results in the figure.

**Comment:** No experiments for this.



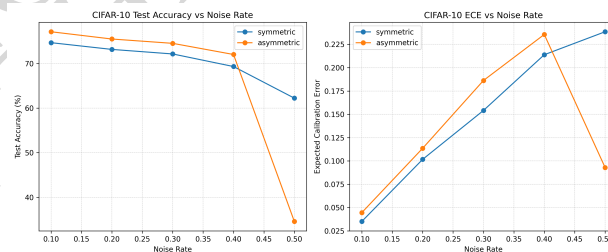
- 216 Jarosław Błasiok and Preetum Nakkiran. Smooth ece: Principled reliability diagrams via kernel  
217 smoothing. *ArXiv*, abs/2309.12236, 2023.  
218
- 219 Aritra Ghosh, Himanshu Kumar, and P. Sastry. Robust loss functions under label noise for deep  
220 neural networks. *ArXiv*, abs/1712.09482, 2017.  
221
- 222 Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.  
223 *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2015.  
224
- 225 Jiayi Huang, Sangwoo Park, and O. Simeone. Calibration-aware bayesian learning. *2023 IEEE 33rd  
226 International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 1–6, 2023.  
227
- 228 Meelis Kull, Miquel Perello-Nieto, Markus Kängsepp, T. S. Filho, Hao Song, and Peter A. Flach.  
229 Beyond temperature scaling: Obtaining well-calibrated multiclass probabilities with dirichlet cal-  
230 ibration. *ArXiv*, abs/1910.12656, 2019.  
231
- 232 Moseli Mots’oehli and kyungim Baek. Gci-vital: Gradual confidence improvement with vision  
233 transformers for active learning on label noise. *ArXiv*, abs/2411.05939, 2024.  
234
- 235 Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy  
236 labels with deep neural networks: A survey. *IEEE Transactions on Neural Networks and Learning  
237 Systems*, 34:8135–8153, 2020.  
238
- 239 Cheng Wang. Calibration in deep learning: A survey of the state-of-the-art. *ArXiv*, abs/2308.01222,  
240 2023.  
241
- 242 Yuan Zhao, Jiayi Chen, and Samet Oymak. On the role of dataset quality and heterogeneity in model  
243 confidence. *ArXiv*, abs/2002.09831, 2020.  
244  
245

## 246 SUPPLEMENTARY MATERIAL

### 247 A ADDITIONAL EXPERIMENTS AND FIGURES

#### 248 A.1 NOISE RATE SENSITIVITY ANALYSIS

249 To provide a deeper understanding of how noise rates affect model performance, we conducted a  
250 noise rate sensitivity analysis on CIFAR-10.



251 Figure 4: CIFAR-10 Test Accuracy vs. Noise Rate for ResNet-18 under symmetric and asymmetric  
252 label noise.

253 Figure 4 shows that as the noise rate increases, test accuracy decreases steadily for both symmetric  
254 and asymmetric noise. The decline is more pronounced under asymmetric noise, reinforcing the  
255 observations made in the main text.

**Comment:**  
This figure is a duplicate of Figure 1, likely because the VLM-based duplication checker overlooked it or the writeup phase failed to account for duplicates.

270  
271  
272  
273  
274  
275  
276  
277

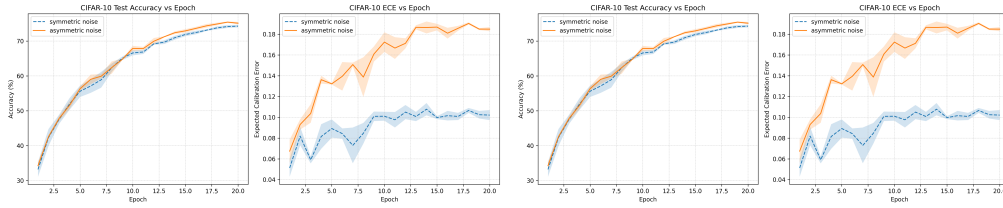


Figure 5: CIFAR-10 Calibration: (Left) Test Accuracy and ECE over epochs; (Right) Aggregated ECE across different noise rates under label noise.

### A.2 CALIBRATION CURVES AND RELIABILITY DIAGRAMS

We also analyzed calibration curves and reliability diagrams to visualize the calibration performance.

Figure 5 illustrates that ECE increases as training progresses, especially under higher noise rates. The reliability diagrams (not shown due to space constraints) further confirm that predictions become overconfident as label noise increases.

### A.3 HYPERPARAMETERS

Table 1 lists the hyperparameters used in our experiments for reproducibility.

Table 1: Hyperparameters used in the experiments.

Parameter	Value
Optimizer	SGD with Momentum
Momentum	0.9
Initial Learning Rate	0.1
Learning Rate Decay	0.1 at epochs 50 and 75
Number of Epochs	100
Batch Size	128
Weight Decay	5e-4

**Comment:** The first two and last two plots are identical. Also these figures are duplicates of the 1st and 4th plots from Figure 2. The y-axis scaling is different, which may explain why the duplication checker missed them.

**Comment:** Weight decay is applied during preliminary experiments only. The experiments use a Cosine Annealing scheduler for the learning rate. The number of epochs is either 20 or 30 instead of 100.

**Comment:** There is no figure for the reliability diagrams, but this time, the writeup phase provided a justification, citing space constraints. This suggests that the system recognizes they are missing.

305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

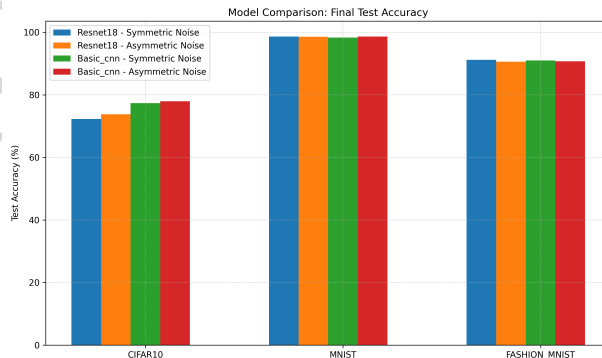


Figure 6: Comparison of Final Test Accuracy between different models under varying noise levels on CIFAR-10.

Figure 6 provides additional insights into how different model architectures perform under label noise, complementing the findings in Section 4.

**Comment:** This figure is a duplicate of Figure 3.

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

A.4 ADDITIONAL DATASETS

We also experimented with SVHN (?), a dataset comprising street view house numbers, to verify the generality of our findings. Results were consistent with previous observations, with label noise adversely affecting calibration metrics.

**Comment:**  
There are no figures for this experiment. The writeup phase should have removed this paragraph.

AUTONOMOUSLY GENERATED BY THE AI SCIENTIST-V2

### C.2.2. AI Scientist Team Review

**Paper Summary** This paper studies the impact of label noise on model calibration using different noise models. More specifically, the paper contrasts symmetric (unstructured label perturbations) and asymmetric (structured label perturbations) noise. The empirical experiments consider standard small-scale vision datasets (i.e. MNIST, Fashion-MNIST and CIFAR-10) and demonstrate that asymmetric noise leads to higher expected calibration error.

#### Strengths

- The research question is of real-world importance and shines light on the impact of noisy labels beyond their effect on prediction accuracy.
- The study design is simple and focuses on a single key factor, i.e. the impact of different noise models (asymmetric noise increasing ECE more than symmetric noise). The considered datasets are appropriate for a workshop submission.
- The impact of the different noise models on the downstream model calibration is robust and consistent across the considered datasets.

#### Weaknesses

- There are multiple instances where the written interpretation of results are not substantially supported by the empirical results presented. E.g. the paragraph interpreting figure 3 refers to ECE measures, which are not displayed in the figures.
- The paper states that it compares different calibration methods, but the paper does not provide any results. The same holds for the mentioned reliability diagrams.
- Furthermore, the supplementary material includes duplicate figures, a missing citation for SVHN and a corresponding missing figure.

#### Scores

- Soundness: 2 fair.  $\Rightarrow$  Interesting research question with potentially simple empirical evaluation setup.
- Presentation: 1 poor.  $\Rightarrow$  Wrong description and duplication of figures. Missing citation and downplaying of related work.
- Contribution: 1 poor.  $\Rightarrow$  While the question considered is important, the displayed results do not provide enough evidence for the conclusions drawn.
- Overall - Workshop: 3/10 (Reject): For instance, a paper with technical flaws, weak evaluation, inadequate reproducibility and incompletely addressed ethical considerations.
- Overall - Conference: 2/10: (Strong reject): For instance, a paper with major technical flaws, and/or poor evaluation, limited impact, poor reproducibility and mostly unaddressed ethical considerations.
- Confidence: 4/5. You are confident in your assessment, but not absolutely certain. It is unlikely, but not impossible, that you did not understand some parts of the submission or that you are unfamiliar with some pieces of related work.

#### Additional Comments

- The biggest flaw of this paper is the mentioning of results that are not substantiated by results. This includes the assessment of various methods tailored to uncertainty calibration, as well as the usage of reliability diagrams. The paper could be substantially improved if these results were added and the selection of displayed figure results was better curated.

- The readability of Figure 2 should be improved by splitting the 6 plots across 2 rows. Furthermore, the related work section appears to dismiss efforts by the scientific community to relate calibration and noisy data.

**Potential Violation of Code of Ethics:** No.

### C.2.3. AI Scientist Team Code Review

#### Temperature scaling

In our review of the paper, we noted that it lacked experiments involving temperature scaling. Upon inspecting the generated code, we found that the AI Scientist had implemented temperature scaling, as can be seen in fig. 8, but never actually used it.

During the paper writing stage, the AI Scientist had access to a set of generated experiment code and its initial plans before generating the code. As a result, it is likely that the paper was influenced by these plans and code, which included temperature scaling, but the AI Scientist failed to realize that the experiments using temperature scaling were never actually conducted.

```
class TemperatureScaling(nn.Module):
    def __init__(self):
        super(TemperatureScaling, self).__init__()
        self.temperature = nn.Parameter(torch.ones(1))

    def forward(self, logits):
        return logits / self.temperature
```

Figure 8 | Temperature scaling implementation.

#### Dataset class

We found that the initial implementation of the dataset class lacked an option for symmetric/asymmetric noise distribution, even though it was part of the initial plan. The AI Scientist recognized this mistake and later implemented the correct version, as shown in fig. 9.

In the main paper, the AI Scientist wrote: “Assymmetric Noise: Labels are flipped to specific incorrect classes based on a predefined confusion matrix, simulating more realistic mislabeling.” The asymmetric noise implementation in the generated code always maps class  $i$  to class  $(i + 1) \% \text{NUM\_CLASSES}$ . While this is a valid approach, it is worth noting that there are other ways to implement asymmetric noise.

```
# Data preparation with noise injection
class NoisyDataset(Dataset):
    def __init__(self, dataset, noise_rate=0.2):
        self.dataset = dataset
        self.noise_rate = noise_rate
        self.noisy_labels = self._inject_noise()

    def _inject_noise(self):
        labels = np.array([y for _, y in self.dataset])
        mask = np.random.rand(len(labels)) < self.noise_rate
        noisy_labels = labels.copy()
        noisy_labels[mask] = np.random.randint(0, NUM_CLASSES, mask.sum())
        return torch.LongTensor(noisy_labels)

    def __getitem__(self, index):
        image, _ = self.dataset[index]
        return image, self.noisy_labels[index]

    def __len__(self):
        return len(self.dataset)

class NoisyDataset(Dataset):
    def __init__(self, dataset, noise_type="symmetric", noise_rate=0.2):
        self.dataset = dataset
        self.noise_rate = noise_rate
        self.noise_type = noise_type
        self.noisy_labels = self._inject_noise()

    def _inject_noise(self):
        labels = np.array([y for _, y in self.dataset])
        if self.noise_type == "symmetric":
            mask = np.random.rand(len(labels)) < self.noise_rate
            noisy_labels = labels.copy()
            noisy_labels[mask] = np.random.randint(0, NUM_CLASSES, mask.sum())
        else: # asymmetric noise
            noisy_labels = labels.copy()
            for i in range(NUM_CLASSES):
                mask = ((labels == i) & (np.random.rand(len(labels)) < self.noise_rate))
                noisy_labels[mask] = (i + 1) % NUM_CLASSES
            return torch.LongTensor(noisy_labels)

    def __getitem__(self, index):
        image, _ = self.dataset[index]
        return image, self.noisy_labels[index]

    def __len__(self):
        return len(self.dataset)
```

Figure 9 | Noisy dataset class implementation.

#### Evaluation function

The evaluation function used to compute the Expected Calibration Error is shown in fig. 10. We manually created test cases and used the MulticlassCalibrationError function with  $\text{norm}='l1'$  from torchmetrics as the ground truth. Since the MulticlassCalibrationError function expects probability inputs, we omitted the softmax operation in the first line to align with the implementation details.

After this adjustment, we confirmed that both functions produce the same results, apart from minor numerical differences.

```
def compute_ece(logits, labels, n_bins=15):
    softmaxes = softmax(logits, dim=1)
    confidences, predictions = torch.max(softmaxes, 1)
    accuracies = predictions.eq(labels)

    bin_boundaries = torch.linspace(0, 1, n_bins + 1)
    bin_lowers = bin_boundaries[:-1]
    bin_uppers = bin_boundaries[1:]

    ece = torch.zeros(1, device=logits.device)
    for bin_lower, bin_upper in zip(bin_lowers, bin_uppers):
        in_bin = confidences.gt(bin_lower.item()) * confidences.le(bin_upper.item())
        prop_in_bin = in_bin.float().mean()
        if prop_in_bin.item() > 0:
            accuracy_in_bin = accuracies[in_bin].float().mean()
            avg_confidence_in_bin = confidences[in_bin].mean()
            ece += torch.abs(avg_confidence_in_bin - accuracy_in_bin) * prop_in_bin
    return ece.item()
```

Figure 10 | The implementation for Expected Calibration Error.

### C.2.4. Workshop Reviews

**Reviewer #1:** This work explores the impact of label noise on model calibration, demonstrating that label noise degrades calibration performance.

The research question is intriguing; however, the experimental analysis appears somewhat unclear. The underlying mechanism explaining how the experimental results support the claimed statement is not well articulated. Specifically, in the abstract, the authors state that "label noise leads to overconfident and miscalibrated predictions, undermining the reliability of uncertainty estimates," yet I struggle to see a clear connection between this claim and the content in the main body.

Additionally, the experimental setup raises some concerns. To thoroughly assess the impact of label noise on model calibration, a more refined approach to introducing label noise should be considered. Moreover, incorporating a broader range of evaluation metrics would help strengthen the conclusions.

Furthermore, the images in the paper are difficult to interpret, and some citations appear to be missing. The referenced papers are also kind of old, which could weaken the soundness of related work.

Rating: 3: Clear rejection

Award: No Award

Confidence: 5: The reviewer is absolutely certain that the evaluation is correct and very familiar with the relevant literature

**Reviewer #2:** Official Review for Submission41

This paper is not finished, there are missing references indicated by (?) there are unlinked references (eg L097), the figures are unreadable (eg Fig 2). I feel this paper is in a late draft status and not review ready.

Remaining comment

It is unclear how Noise Mitigation Techniques, or Calibration Improvements (like temperature scaling) are taken into account in the study. Or how they affect performance after label noise. It is stated that temperature scaling is used, but its effect is not made clear.

Rating: 3: Clear rejection

Award: No Award

Confidence: 4: The reviewer is confident but not absolutely certain that the evaluation is correct



### C.3. Real-world Challenges in Pest Detection using Deep Learning: an Investigation into Failures and Solutions

#### C.3.1. THE AI SCIENTIST-v2 Idea

##### Idea

```

>Name": "real_world_pest_detection",
>Title": "Real-World Challenges in Pest Detection Using Deep Learning: An Investigation into Failures and Solutions",
'Short Hypothesis': "Deep learning models for pest detection often fail to generalize in real-world agricultural settings due to data quality issues, environmental variability, and model limitations. Investigating these failures can lead to more robust solutions.",
'Related Work': "Several studies, such as those by Agarwal et al. (2023) and Dong et al. (2024), have explored deep learning for pest detection in agriculture. These studies generally report high accuracy in controlled settings but often do not address real-world deployment challenges. Our proposal distinguishes itself by focusing on the negative outcomes and the underlying reasons behind these failures.",
'Abstract': "Accurate pest detection is vital for protecting crops and ensuring food security. While deep learning models have shown promise in controlled environments, their performance often degrades in real-world applications. This proposal aims to investigate the reasons behind these failures. We hypothesize that data quality issues, environmental variability, and model limitations are significant factors. By conducting a series of experiments, we will explore these challenges in depth and propose robust solutions to improve the generalizability of deep learning models for pest detection. Our research will provide valuable insights for the agricultural community and contribute to the development of more reliable AI tools for precision farming.",
'Experiments': [
  "1. **Data Quality Analysis**: Collect a diverse dataset of pest images from different agricultural environments and analyze its quality. Identify common issues such as label noise, class imbalance, and distribution shift.",
  "2. **Model Robustness Testing**: Train state-of-the-art deep learning models (e.g., YOLOv8, EfficientNetB3) on the collected dataset and evaluate their performance in controlled vs. real-world settings. Metrics: Mean Average Precision (mAP), F1 Score.",
  "3. **Environmental Variability Study**: Evaluate model performance under different environmental conditions (e.g., lighting, weather). Identify which conditions most significantly impact model accuracy.",
  "4. **Failure Mode Analysis**: Conduct a detailed analysis of misclassifications to identify common patterns and potential causes (e.g., feature overlap between pests and background).",
  "5. **Improvement Strategies**: Implement and test various strategies to mitigate identified challenges, such as data augmentation, domain adaptation, and model ensembling. Evaluate their effectiveness in improving model robustness."
],
'Risk Factors and Limitations': "Potential risks include the availability and quality of real-world data, the computational demands of training and testing multiple deep learning models, and the generalizability of the findings to different types of pests and crops. Additionally, environmental factors may introduce variability that is challenging to control.",
'Code': "from datasets import ..."

```

**Link to more material:** <https://github.com/SakanaAI/AI-Scientist-ICLR2025-Workshop-Experiment/tree/master/pest-detection>.

# REAL-WORLD CHALLENGES IN PEST DETECTION USING DEEP LEARNING: AN INVESTIGATION INTO FAILURES AND SOLUTIONS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Deep learning models have shown significant promise in pest detection tasks within controlled environments, but their performance often degrades when deployed in real-world agricultural settings. This study investigates the challenges hindering the generalization of these models, focusing on data quality issues, environmental variability, and inherent model limitations. Through extensive experiments, including learning rate optimization and multi-dataset training, we demonstrate that while lower learning rates can enhance generalization, models still struggle with robustness to environmental changes. Our findings highlight critical pitfalls in deploying deep learning models for pest detection and offer insights into potential solutions for improving their real-world applicability.

**Comment:**  
The rest of the paper primarily focuses on environmental variability and noise in images.

## 1 INTRODUCTION

Accurate pest detection is crucial for protecting crops and ensuring global food security. Deep learning models have emerged as powerful tools for automating pest detection tasks, achieving high accuracy in controlled environments. However, their performance often degrades significantly when deployed in real-world agricultural settings. This gap between controlled experiments and real-world applications poses a serious challenge for precision agriculture and highlights the need for robust, generalizable models. Understanding and addressing the reasons behind these performance drops is essential for advancing AI in agriculture.

**Comment:**  
Both sentences require citations to be substantiated.

**Comment:**  
This might be setting an unnecessarily high bar.

In this work, we investigate the factors contributing to the failures of deep learning models in real-world pest detection scenarios. We hypothesize that issues such as data quality, environmental variability, and inherent model limitations play significant roles in hindering model generalization. Through a series of experiments, we explore these challenges in depth. Our findings reveal that while optimizing hyperparameters, such as the learning rate, can lead to improved validation accuracy, deep learning models still struggle to maintain robustness under environmental changes. Moreover, multi-dataset training and domain adaptation techniques, aimed at enhancing generalization across different datasets, present their own set of challenges, including increased computational demands and inconsistent performance gains.

**Comment:**  
Is it really in depth?

By presenting these negative and inconclusive results, we aim to highlight the real-world pitfalls and challenges in deploying deep learning models for pest detection. Our research provides valuable insights for the agricultural and machine learning communities, contributing to the development of more reliable AI tools for precision farming.

**Comment:**  
The domain adaptation experiment in Section 5.2 appears to highlight the challenges of adapting ImageNet-trained models to other vision datasets when environmental noise is present, but this experiment isn't directly related to the main focus of the paper.

## 2 RELATED WORK

Deep learning has been widely applied in agricultural contexts for tasks such as pest and disease detection, showing high accuracy in controlled settings (Mustakim et al., 2024; Kumar et al., 2022). Li et al. (2023b) highlighted the limitations of traditional deep learning methods in practical applications, noting issues such as overfitting and sensitivity to environmental variations. Several studies have explored methods to improve model robustness and generalization. Data augmentation techniques have been employed to enhance dataset diversity and reduce overfitting (Abdulkareem

et al., 2024). Domain adaptation strategies have been proposed to address domain shifts and improve performance in new environments (Prasad & Agniraj, 2024; Li et al., 2023a). However, these approaches often do not fully address the challenges faced in real-world deployment.

Reviews like Teixeira et al. (2023) and Hu (2023) have identified gaps in current research, emphasizing the need for models that generalize well to diverse, real-world conditions. Additionally, Amir et al. (2024) discussed the limitations of deep learning models when encountering out-of-distribution inputs, underscoring the importance of verifying model generalization. Our work distinguishes itself by focusing on the failures and limitations of deep learning models in real-world pest detection scenarios, providing an in-depth investigation into the underlying causes and proposing insights for improvement.

### 3 METHODOLOGY

We employed deep learning models for pest detection, focusing on evaluating their performance and robustness in real-world agricultural settings. We utilized the ResNet-18 architecture (?), pretrained on ImageNet, and fine-tuned it on the Crop Pest and Disease dataset, which includes 22 classes of pests and diseases collected from local farms. To investigate the challenges, we designed experiments to assess the impact of learning rates on model performance. We hypothesized that optimizing the learning rate could improve generalization. We also implemented data augmentation techniques to simulate environmental variability, such as brightness and contrast changes, Gaussian blur, and random affine transformations, to evaluate the models' robustness. Additionally, we explored multi-dataset training using datasets such as EuroSAT (?), MedMNIST (?), and CIFAR-10 (?) to assess the potential of domain adaptation and transfer learning in improving model generalization across different agricultural domains.

### 4 EXPERIMENTAL SETUP

The Crop Pest and Disease dataset comprises 25,126 images across 22 classes of pests and diseases affecting crops such as cashew, cassava, maize, and tomato. We split the dataset into training (70%), validation (15%), and testing (15%) sets. For the baseline experiments, we conducted a grid search to optimize the learning rate, evaluating values of  $\{1e^{-4}, 5e^{-4}, 1e^{-3}, 5e^{-3}, 1e^{-2}\}$ . We trained the ResNet-18 model for 10 epochs for each learning rate, using a batch size of 32 and the Adam optimizer. To simulate challenging environmental conditions, we applied data augmentations during testing, including brightness and contrast adjustments, Gaussian blur, and random affine transformations. We introduced the Environmental Robustness Score (ERS), calculated as the ratio of model accuracy under challenging conditions to that under normal conditions, to quantify robustness.

In the research experiments, we trained models on additional datasets—EuroSAT, MedMNIST, and CIFAR-10—to investigate the effects of multi-dataset training on model generalization. We used similar training settings and evaluated models using accuracy, loss, and ERS.

## 5 EXPERIMENTS AND RESULTS

### 5.1 IMPACT OF LEARNING RATE ON MODEL PERFORMANCE

To evaluate the impact of learning rates on model performance, we trained the ResNet-18 model on the Crop Pest and Disease dataset using different learning rates. Figure 1 illustrates the aggregated accuracy, loss, and ERS across different learning rates.

As shown in Figure 1, lower learning rates ( $1e^{-4}$  and  $5e^{-4}$ ) result in smoother convergence of training and validation accuracy, and a steady decrease in training loss. The ERS remains more stable for these learning rates, suggesting enhanced robustness to environmental variability. In contrast, higher learning rates lead to overfitting and unstable loss patterns, with significant fluctuations in ERS scores. These results indicate that optimizing hyperparameters like learning rate is crucial for improving model generalization and robustness in real-world settings.

**Comment:**  
Cite. Also would be nice to show some example images in the appendix.

**Comment:**  
EuroSAT might make sense but using MedMNIST and CIFAR-10 to "assess the potential of domain adaptation...across different agricultural domains" is a stretch b/c these two are clearly not from the agricultural domain.

**Comment:**  
The figure seems to suggest the lower lr group overfits more

**Comment:**  
The citation key 'he2016deep' exists in the latex file but not in references.bib.

**Comment:** All three datasets citation key exist in the latex file but not in the bib file.

**Comment:**  
used a subset (around 2500 images)

**Comment:**  
only in this specific condition and definition of environmental robustness

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141

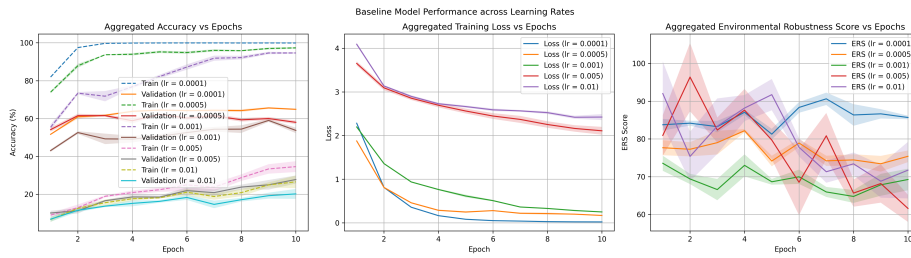


Figure 1: Baseline model performance across learning rates. Aggregated training and validation accuracy, training loss, and Environmental Robustness Score (ERS) over epochs for different learning rates. Lower learning rates yield higher validation accuracy and more stable ERS scores, indicating better generalization and robustness.

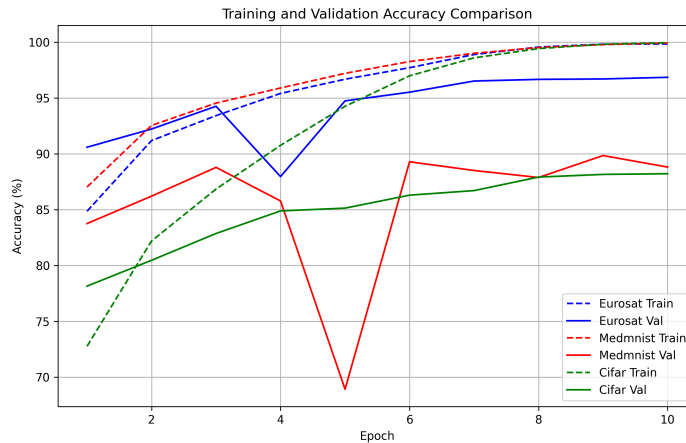


Figure 2: Comparison of training and validation accuracy across different datasets. Models trained on EuroSAT and CIFAR-10 exhibit stable and high accuracy, whereas the model trained on MedMNIST shows erratic accuracy patterns, indicating challenges in generalization due to domain discrepancies.

## 5.2 CHALLENGES IN MULTI-DATASET TRAINING

To investigate model generalization across different domains, we trained the ResNet-18 model on additional datasets: EuroSAT, MedMNIST, and CIFAR-10. Figure 2 presents the comparison of training and validation accuracy across these datasets.

From Figure 2, the models trained on EuroSAT and CIFAR-10 achieve high and stable training and validation accuracy over epochs, suggesting effective learning and better generalization. In contrast, the MedMNIST model displays fluctuating accuracy, highlighting difficulties in adapting to the pest detection task. This suggests that significant domain shifts can negatively impact learning, leading to overfitting and decreased robustness.

To further examine the robustness of these models, we analyzed the ERS across epochs, as shown in Figure 3.

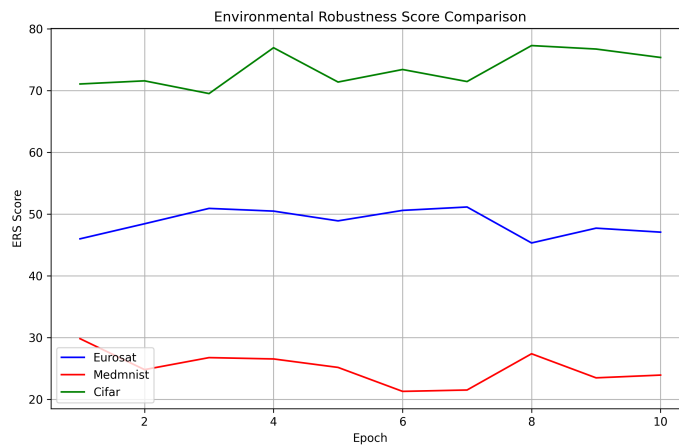
Figure 3 illustrates that the models trained on EuroSAT and CIFAR-10 maintain higher ERS scores across epochs, suggesting robustness to environmental augmentations applied during testing. The MedMNIST model's low ERS scores indicate vulnerability to such changes, underscoring the challenges posed by domain differences. These findings highlight the varying impact of dataset characteristics on model generalization and robustness. While multi-dataset training and domain adapta-

**Comment:** Should explain that this is more about generalization from ImageNet to EuroSAT, MedMNIST, or CIFAR.

149  
150  
151  
152

**Comment:** a stretched argument. Difficulty of adopting ResNet-18 to MedMNIST doesn't immediately suggest it's also difficult to adapting to the pest detection task. At least it needs more explanation.

**Comment:** Should refer to Appendix B, as it provides more detailed information on the noise used



**Comment:**  
The legend can be improved

Figure 3: Environmental Robustness Score (ERS) comparison across different datasets. The EuroSAT and CIFAR-10 models maintain higher and more stable ERS scores, indicating better robustness to environmental variability. The MedMNIST model shows low and unstable ERS scores, reflecting sensitivity to environmental changes.

tion can offer potential improvements, they also present challenges, such as increased computational demands and inconsistent performance gains, which must be carefully managed.

## 6 DISCUSSION

Our experiments reveal significant challenges in deploying deep learning models for pest detection in real-world agricultural settings. Optimizing hyperparameters like learning rates enhances model generalization and robustness to some extent, as evidenced by the improved performance and stable ERS scores at lower learning rates. However, models still struggle with environmental variability, indicating that hyperparameter optimization alone is insufficient to achieve robust real-world performance.

The exploration of multi-dataset training provides valuable insights into domain adaptation challenges. The varying performance across datasets underscores the importance of dataset selection and the potential pitfalls of naively combining datasets with different characteristics. The model's poor performance on MedMNIST suggests that significant domain shifts can negatively impact learning, leading to overfitting and decreased robustness.

These findings emphasize the need for specialized strategies to address data quality issues and environmental variability. Data augmentation techniques that more accurately reflect real-world conditions, robust training methods, and domain-specific model adaptations may be necessary to improve model performance in practical applications.

## 7 CONCLUSION

Our study highlights critical pitfalls in deploying deep learning models for pest detection in real-world agricultural settings. While optimizing hyperparameters like learning rates can enhance model generalization and robustness, challenges remain due to environmental variability and domain discrepancies. Multi-dataset training introduces additional complexities, and its benefits depend on the compatibility of the datasets involved. Future work should focus on developing advanced techniques tailored to real-world conditions, such as improved data augmentation strategies that mimic environmental changes, robust training methods that enhance model resilience, and architectures designed for adaptability. By addressing these challenges, we can move closer to deploying reliable AI tools in precision agriculture that are resilient to real-world variability.

**Comment:**  
This can be misleading, since multi-dataset training could mean a model trained on multiple datasets, but here multiple models are trained, where each model is trained on a single dataset

**Comment:**  
This only makes sense if "combining" refers to ImageNet pretraining followed by fine-tuning on a different dataset, but it usually gives the impression that the model is trained on multiple datasets.

## REFERENCES

- Ismael M. Abdulkareem, Faris K. Al-Shammri, Noor Aldeen A. Khalid, and Natiq A. Omran. Proposed approach for object detection and recognition by deep learning models using data augmentation. *Int. J. Online Biomed. Eng.*, 20:31–43, 2024.
- Guy Amir, Osher Maayan, Tom Zelazny, Guy Katz, and Michael Schapira. Verifying the generalization of deep learning to out-of-distribution domains. *ArXiv*, abs/2406.02024, 2024.
- Jiangfeng Hu. Application of deep learning in smart agriculture research. *Applied and Computational Engineering*, 2023.
- Raj Kumar, Dinesh Singh, A. Chug, and A. Singh. Evaluation of deep learning based resnet-50 for plant disease classification with stability analysis. In *International Conference Intelligent Computing and Control Systems*, pp. 1280–1287, 2022.
- A. Li, Elisa Bertino, Rih-Teng Wu, and Ting Wu. Building manufacturing deep learning models with minimal and imbalanced training data using domain adaptation and data augmentation. *2023 IEEE International Conference on Industrial Technology (ICIT)*, pp. 1–8, 2023a.
- Manzhou Li, Siyu Cheng, Jingyi Cui, Changxiang Li, Zeyu Li, Chang Zhou, and Chunli Lv. High-performance plant pest and disease detection based on model ensemble with inception module and cluster algorithm. *Plants*, 12, 2023b.
- M. Mustakim, Aditya Rezky Pratama, Imam Ahmad, Teguh Arifianto, Kelik Sussolaikah, and Sepriano Sepriano. Image classification of corn leaf diseases using cnn architecture resnet-50 and data augmentation. In *2024 International Conference on Decision Aid Sciences and Applications (DASA)*, pp. 1–6, 2024.
- Pulicherla Siva Prasad and Senthilrajan Agniraj. Cross-domain adaptation techniques for robust plant disease detection: A dann-coral hybrid approach. *International Journal of Experimental Research and Review*, 2024.
- A. Teixeira, José Ribeiro, R. Morais, J. Sousa, and António Cunha. A systematic review on automatic insect detection using deep learning. *Agriculture*, 2023.

## SUPPLEMENTARY MATERIAL

### A ADDITIONAL FIGURES AND DETAILED RESULTS

We provide additional figures and detailed results to supplement the main text. These figures offer deeper insights into the models’ behaviors under different experimental conditions.

Figure 4 shows that the EuroSAT and CIFAR-10 models show consistent decreases in training loss, reflecting effective learning. The MedMNIST model’s erratic loss suggests that the model struggles to minimize the loss function, possibly due to significant differences between medical images and agricultural pest images.

### B IMPLEMENTATION DETAILS

All models were implemented using PyTorch 1.9.0. The ResNet-18 architecture was initialized with ImageNet pretrained weights. For optimization, we used the Adam optimizer with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , and a weight decay of  $1e^{-4}$ . No learning rate schedules or gradient clipping were applied.

Data augmentations for simulating challenging conditions included `ColorJitter` with brightness and contrast factors of 0.5, `GaussianBlur` with a kernel size of 3, and `RandomAffine` transformations with degrees up to 15 and translation up to 10%. These augmentations were applied during testing to evaluate the Environmental Robustness Score (ERS).

**Comment:**  
This should be ImageNet images

**Comment:**  
weight decay is 0.01 instead of  $1e^{-4}$ .

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300

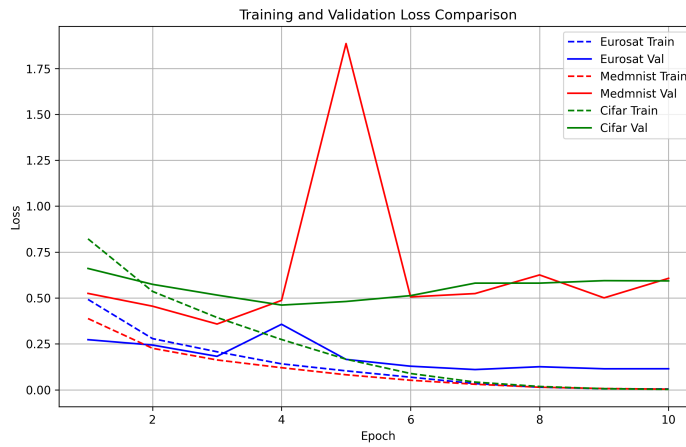


Figure 4: Comparison of training loss across different datasets. Models trained on EuroSAT and CIFAR-10 datasets demonstrate a steady decrease in loss, while the model trained on MedMNIST exhibits erratic loss curves, indicating instability during training due to domain mismatch.

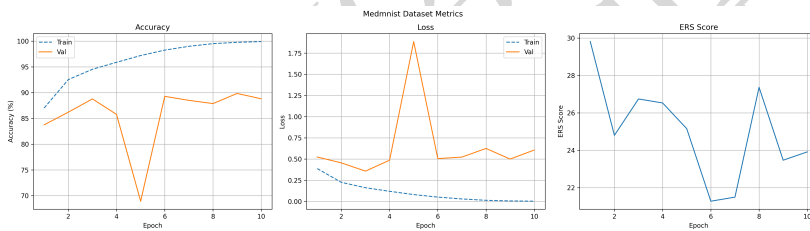


Figure 5: Performance metrics for the model trained on MedMNIST dataset. The erratic behavior in accuracy and loss indicates challenges in model convergence and generalization when applying the MedMNIST dataset to pest detection tasks.

**Comment:**  
This statement is incorrect again. It should say something like "...applying/using the MedMNIST dataset with the ImageNet-pretrained model." Also, the figure isn't mentioned in the main text.

The Environmental Robustness Score (ERS) is defined as:

$$ERS = \frac{\text{Accuracy under challenging conditions}}{\text{Accuracy under normal conditions}} \quad (1)$$

This metric quantifies the model's robustness to environmental changes by comparing its performance under augmented test sets to that under standard conditions.

The additional datasets used for multi-dataset training were:

- **EuroSAT:** A dataset consisting of 27,000 labeled Sentinel-2 satellite images covering 10 classes (?).
- **MedMNIST:** A collection of lightweight medical image datasets covering various tasks (?).
- **CIFAR-10:** A well-known dataset consisting of 60,000 32x32 color images in 10 classes (?).

For the multi-dataset training, we used a batch size of 64 to accommodate the increased data volume. Training was conducted for 30 epochs, and early stopping was applied if validation loss did not decrease for 5 consecutive epochs.

314  
315  
316  
317  
318  
319  
320  
321  
322  
323

### C.3.2. AI Scientist Team Review

**Paper Summary** This paper studies the application of Deep Learning models for a real-world application to pest prediction. It introduces an Environmental Robustness Score that leverages various data augmentation techniques, mimicking environmental factors affecting data collection. It compares various learning rates and compares the impact of out-of-distribution testing settings across non-pest vision datasets.

#### Strengths

- The paper fits the ICBNB workshop topic especially well. It discusses a real-world application of Deep Learning methods to pest prediction.
- Understanding the differential impact of training and out-of-distribution data augmentation technique settings across datasets is interesting.

#### Weaknesses

- The paper refers to domain adaptation being studied multiple times. The experiments, on the other hand, only investigate the usage of data augmentation methods (such as lighting, blurring, and contrast manipulation). Furthermore, studying the impact of the learning rate on generalization is fairly trivial.
- It is hard to motivate that the Eurosat, Medmnist, and CIFAR-10 results are related to the pest prediction problem. Why should a result on these datasets transfer to pest prediction?
- Some of the statements regarding multi-dataset training are misleading. There are no results in the paper that result from such a training setup. Instead, multiple models are trained on individual datasets.

#### Scores

- Soundness: 2 fair.  $\Rightarrow$  Interesting research question with potentially simple empirical evaluation setup.
- Presentation: 1 poor.  $\Rightarrow$  Wrong description and duplication of figures. Missing citation and downplaying of related work.
- Contribution: 1 poor.  $\Rightarrow$  While the question considered is important, the displayed results do not provide enough evidence for the conclusions drawn.
- Overall - Workshop: 3/10 (Reject): For instance, a paper with technical flaws, weak evaluation, inadequate reproducibility and incompletely addressed ethical considerations.
- Overall - Conference: 2/10: (Strong reject): For instance, a paper with major technical flaws, and/or poor evaluation, limited impact, poor reproducibility and mostly unaddressed ethical considerations.
- Confidence: 4/5. You are confident in your assessment, but not absolutely certain. It is unlikely, but not impossible, that you did not understand some parts of the submission or that you are unfamiliar with some pieces of related work.

#### Additional Comments

- The presentation of the results needs significant improvement. There are multiple missing citations (?), and the interpretation of the results can be misleading. This includes the conclusions with regard to the impact of a lower learning rate on overfitting or naming the multi-model-single-dataset experiment “mult-dataset”.

**Potential Violation of Code of Ethics:** No.



### C.3.3. AI Scientist Team Code Review

#### Domain Adaptation and Multi-dataset training

The paper seems to describe the “Domain Adaptation” experiment as primarily focused on transferring ImageNet-pretrained models to other vision datasets. After reviewing the code, we found attempts to implement a domain adaptation technique by training a separate classifier to distinguish different domains, but these attempts were unsuccessful. In the end, the AI Scientist opted for an implementation that does not include this domain adaptation technique.

Moreover, in the code where this domain adaptation technique was implemented, multi-dataset training was correctly performed as well—training a single model on all three datasets with domain discriminator loss, as shown in fig. 11. Had this code run successfully, the AI Scientist would likely have chosen it over the one ultimately selected, which lacked proper multi-dataset training but ran without errors.

```

class DomainDiscriminator(nn.Module):
    def __init__(self):
        super().__init__()
        self.fc1 = nn.Linear(FEATURE_DIM, 256)
        self.fc2 = nn.Linear(256, 3)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        return self.fc2(x)

```

```

# Training loop
for epoch in range(NUM_EPOCHS):
    feature_extractor.train()
    domain_discriminator.train()
    classifier.train()

for dataset_name, (train_dataset, test_dataset) in datasets.items():
    train_loader = DataLoader(train_dataset, batch_size=BATCH_SIZE, shuffle=True)

    for batch in train_loader:
        # Move batch to device
        images = batch["image"].to(device)
        labels = batch["label"].to(device)

        # Feature extraction
        features = feature_extractor(images)

        # Classification loss
        clf_outputs = classifier(features)
        clf_loss = F.cross_entropy(clf_outputs, labels)

        # Domain adversarial loss
        domain_outputs = domain_discriminator(features)
        domain_labels = torch.zeros(images.size(0), dtype=torch.long).to(device)
        domain_loss = F.cross_entropy(domain_outputs, domain_labels)

        # Total loss
        total_loss = clf_loss - 0.1 * domain_loss

        # Optimization
        fe_optimizer.zero_grad()
        clf_optimizer.zero_grad()
        disc_optimizer.zero_grad()
        total_loss.backward()
        fe_optimizer.step()
        clf_optimizer.step()
        disc_optimizer.step()

```

Figure 11 | Domain discriminator and multi-dataset training loop.

#### Environmental noise implementation

The paper states, “To simulate challenging environmental conditions, we applied data augmentations during testing, including brightness and contrast adjustments, Gaussian blur, and random affine transformations.” This is confirmed in the code, as shown in fig. 12.

The calculation of the Environmental Robustness Score—a metric introduced by the AI Scientist and defined as “the ratio of model accuracy under challenging conditions to that under normal conditions, to quantify robustness”—matches the description in the paper, as shown in fig. 13.

```

# Transforms
base_transform = T.Compose(
    [
        T.Resize((64, 64)),
        T.ToTensor(),
        T.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
    ]
)

challenging_transform = T.Compose(
    [
        T.Resize((64, 64)),
        T.ColorJitter(brightness=0.5, contrast=0.5),
        T.RandomAffine(degrees=15, translate=(0.1, 0.1)),
        T.GaussianBlur(kernel_size=3),
        T.ToTensor(),
        T.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
    ]
)

```

Figure 12 | Environment noise simulation implementation.

```

def calculate_ers(model, normal_loader, challenging_loader):
    model.eval()
    with torch.no_grad():
        # Normal conditions
        correct_normal = 0
        total_normal = 0
        for images, labels in normal_loader:
            images, labels = images.to(device), labels.to(device)
            outputs = model(images)
            _, predicted = torch.max(outputs.data, 1)
            total_normal += labels.size(0)
            correct_normal += (predicted == labels).sum().item()
        normal_acc = correct_normal / total_normal

        # Challenging conditions
        correct_challenging = 0
        total_challenging = 0
        for images, labels in challenging_loader:
            images, labels = images.to(device), labels.to(device)
            outputs = model(images)
            _, predicted = torch.max(outputs.data, 1)
            total_challenging += labels.size(0)
            correct_challenging += (predicted == labels).sum().item()
        challenging_acc = correct_challenging / total_challenging

    ers = (challenging_acc / normal_acc) * 100
    return ers

```

Figure 13 | Environmental Robustness Score calculation.

### C.3.4. Workshop Reviews

#### Reviewer #1: Review of Real-World Challenges in Pest Detection Using Deep Learning: An Investigation into Failures and Solutions.

##### Summary

This paper studies deep learning methods in pest detection applications. It highlights the need for more research and attempts to perform first experiments in this area.

##### Results

The paper performs experiments on an image classifier, a ResNet-18 trained on ImageNet, fine-tuned on the Crop Pest and Disease dataset. It uses various data augmentations to emulate the real-world conditions and further trains the model using "multi-dataset training". The model performance is measured in accuracy, loss, and Environmental Robustness Score (ERS). The first experiment investigates the effect tuning the learning rate has on training. The paper claims that a lower learning rate leads to a higher generalisation. The second experiment investigates

the model's generalisation across different datasets. The paper claims that training the model on EuroSAT and CIFAR-10 datasets leads to better generalisation.

#### Strengths

The paper's background, motivation, and related work are well written. The motivation to study generalisation of deep learning methods in real-world agricultural applications is good.

#### Weaknesses

- The experiments are unmotivated and unclear.
  - It is unclear how the choice of augmentation methods are related to "real-world environmental variability."
  - The choice of ERS is unmotivated and no intuition on the score is given in the paper.
  - The learning rate experiment conclusion that the model's generalization and robustness to real-world setting seems misleading since only 5 learning rates are used and the models are only trained for 10 epochs. Furthermore, the model was not tested in a real-world deployment.
- It is unclear what the paper means by "multi-dataset training" especially since the datasets have a different number of classes. Thus, the results of this experiment are unclear.
- The paper claims to have studied "deploying deep learning models for pest detection in real-world agricultural settings", however, the paper does not test the trained model in a real-world setting. Thus, the paper's conclusions are misleading.

#### General Remarks

- In the introduction, the difference between "controlled environment" and "real-world agricultural settings" should be explained since the experiments are not performed in a real-world deployment.
- Plots are small and difficult to read. Increasing the font size would help as well.
- In Figure 1, it would be helpful if the Train and Validation lines for the same learning rate were the same colour.
- In Figure 1, it is unclear whether the ERS scores are evaluated on the train or validation dataset.
- Unclear why the results in Figure 4 are pushed to the appendix and not combined with Figure 2 similar to Figure 1.
- References to the datasets are missing, e.g., the Crop Pest and Disease dataset.

#### Concluding Remarks

Overall, the paper addresses an interesting real-world problem. However, the lack of detail in the experiment section limits the strength of the conclusions, as the experimental results are not sufficiently supported by evidence. In its current state, it is of the reviewer's opinion that the paper does not meet the standard for publication. The authors are encouraged to further develop this research and consider deploying the model in a real-world setting to strengthen the validity of their findings.

Rating: 3: Clear rejection

Award: No Award

Confidence: 4: The reviewer is confident but not absolutely certain that the evaluation is correct

### Reviewer #2: Review "REAL-WORLD CHALLENGES IN PEST DETECTION USING DEEP LEARNING: AN INVESTIGATION INTO FAILURES AND SOLUTIONS"

**Summary:** The authors investigate deep learning models in the context of pest detection. They state that common deep learning models work well in theoretical settings but struggle to generalize when being exposed to environmental changes. In addition, they offer potential approaches to address these issues and increase robustness of deep learning models in pest detection.

**Scientific Rigor and Transparency:** Yes, the authors conducted several experiments to underline their findings.

**Novelty and Significance:** Yes, the paper highlights the weaknesses of deep learning models in pest detection by analyzing how multi-dataset training and hyperparameter tuning affect their performance and ability to generalize.

**Clarity of Writing:** Yes, the paper is generally well-structured and written in a nice way. However, the paper could still improve on clarity by adding the missing references marked with a (?) in Section 3.

**Alignment with Workshop Topics:** Yes, the paper aligns with the theme of the workshop.

Additional Comments: The submitted paper provides insights into the weaknesses of deep learning models in real-world pest detection scenarios. It proposes strategies to mitigate these issues through hyperparameter tuning and multi-dataset training. While these methods can enhance model performance in practical applications, challenges persist due to environmental variability and domain discrepancies. I recommend that the authors include additional references in the field of pest detection, such as "Crop Pest Recognition in Real Agricultural Environments Using Convolutional Neural Networks with a Parallel Attention Mechanism" by Zhao et al., to offer a more comprehensive perspective on the topic.

Rating: 7: Good paper, accept

Award: No Award

Confidence: 3: The reviewer is fairly confident that the evaluation is correct

### Reviewer #3: Critical Review of Real-World Challenges in Pest Detection Using Deep Learning: Methodological and Theoretical Considerations

The presented paper discusses challenges in pest detection based on digital images using the ResNet-18 model. The authors discuss experiments to evaluate the variability of classification performance based on simulated environmental changes. This topic is relevant, given major challenges such as biodiversity loss. Furthermore, the importance of interdisciplinary research (in this case, data science and biology) will increase, and such studies will help accelerate the use of machine learning in life sciences. However, I found shortcomings in this study, which I summarize in the text below.

The introduction provides a good overview of why this work is important. However, the technical motivation is not clear. A clear motivation based on the theoretical aspects of 'generalization' (see [1,2]), as well as a clear statement including literature on challenges in 'AI' in agriculture, would have been necessary.

The methodology section should refer to the appendix for more details (there are important details in the appendix). Furthermore, dataset details (e.g., example images, how many disease-related images are there?) are missing. The hypothesis concerning the learning rate, as well as augmentation to simulate real environmental variability, is not well motivated. I believe this harsh simulation of real-life dynamics should have been introduced in the abstract and introduction (it would still be an interesting study!). The motivation and derivation of the ERS are missing (is it an ad hoc approach?). The metric is prone to over/underestimating robustness due to unbalanced datasets (see the equation, and using the definition of accuracy, the size of the datasets influences the fraction). Based on the missing training/test/validation details above, it is unclear whether bias is introduced. Furthermore, I do not think that such studies must rely on the newest models. However, ResNet-18 is a rather old model, and no justification for selecting this model is given. A comparison to transformer-based architectures would have been interesting.

Finally, there are some language issues and BibTeX errors (see '?'). The figures should be updated to increase readability. Considering my discussion above, I think the results are still interesting. However, I do believe that the presentation must be adapted. I recommend a major revision, including a solid theoretical foundation, a presentation of the evaluation strategy using augmentation throughout the manuscript, and a comparison to recent deep learning models. Furthermore, I recommend switching from augmentation to real datasets or generative models. With these improvements, the impact of this study would be increased significantly.

[1] Wolpert, D.H. (2002). The Supervised Learning No-Free-Lunch Theorems. In: Soft Computing and Industry. Springer, London. [https://doi.org/10.1007/978-1-4471-0123-9\\_3](https://doi.org/10.1007/978-1-4471-0123-9_3) [2] Goldblum, M. et al. (2024). Position: The No Free Lunch Theorem, Kolmogorov Complexity, and the Role of Inductive Biases in Machine Learning. Proceedings of the 41st International Conference on Machine Learning

Rating: 4: Ok but not good enough - rejection

Award: No Award

Confidence: 4: The reviewer is confident but not absolutely certain that the evaluation is correct