# AutoIdeaFlow: from Idea Generation to Paper Writeup and Review

**A Minor Project Report**
Submitted To



**Chhattisgarh Swami Vivekanand Technical University**
**Bhilai, India**
For
Minor Project
of
**Bachelor of Technology (Hons.)**
*in*
**Computer Science & Engineering**
*By*

**Jayant Patel**
**300012721061**
**CB4646**
**7th Sem**
**Artificial Intelligence**

Under the Guidance of
**Dr. Nachiket Tapas**
Assistant Professor
Department of Computer Science & Engineering
**UTD, CSVTU, Bhilai (C.G.)**



**Department of Computer Science & Engineering**
**University Teaching Department**
**Chhattisgarh Swami Vivekanand Technical University**
**Bhilai (C.G.) 491107**

**Session: 2024 – 2025**

# DECLARATION BY THE CANDIDATE

We the undersigned solemnly declare that the Minor project report entitled *"AUTOIDEAFLOW FROM IDEA GENERATION TO PAPER WRITEUP AND REVIEW"* is based our own work carried out during the course of our study under the supervision of *Dr. Nachiket Tapas*.

We assert that the statements made and conclusions drawn are an outcome of the project work. We further declare that to the best of our knowledge and belief that the report does not contain any part of any work which has been submitted for the award of any other degree/diploma/certificate in this University/Deemed university of India or any other country.

**Jayant Patel**
Roll No: 300012721061
Enroll No: CB4646
Semester: 7<sup>th</sup> (CSE)

# CERTIFICATE BY THE SUPERVISOR

This is to certify that the Minor project report entitled *"AUTOIDEAFLOW FROM IDEA GENERATION TO PAPER WRITEUP AND REVIEW"* is a record of project work carried out under my guidance and supervision for the fulfillment of the award of degree of Bachelor of Technology (Hons.) in the faculty of Computer Science & Engineering of Chhattisgarh Swami Vivekananda Technical University, Bhilai (C.G.) India.

To the best of my knowledge and belief the report

  i. Embodies the work of the candidate himself

 ii. Has duly been completed

iii. Fulfills the partial requirement of the ordinance relating to the B.Tech. (Hons) degree of the University

 iv. Is up to the desired standard both in respect of contents and language for being referred to the examiners.

**Dr. Nachiket Tapas**
Assistant Professor
Department of Computer Science &
Engineering, UTD, CSVTU, Bhilai (C.G.)

**Forwared to**
**Chhattisgarh Swami Vivekanand Technical University, Bhilai (C.G.)**

**Dr. J P Patra**
HOD
Dept. of Computer Science &
Engineering UTD, CSVTU, Bhilai
(C.G.)

**Dr. P K Ghosh**
Director
UTD, CSVTU, Bhilai (C.G.)

II

# CERTIFICATE BY THE EXAMINERS

The project report entitled **_"AUTOIDEAFLOW FROM IDEA GENERATION TO PAPER WRITEUP AND REVIEW"_** has been examined by the undersigned as a part of the examination of Bachelor of Technology (Hons.)  in the faculty of Computer Science & Engineering of Chhattisgarh Swami Vivekanand Technical University, Bhilai.

———————————————                                 ———————————————
     **Internal Examiner**                                                        **External Examiner**
        **Date:**                                                                            **Date:**

# ACKNOWLEDGEMENT

Working for this project has been a great experience for us. There were moments of anxiety, when we could not solve a problem for the several days. But we have enjoyed every bit of process and are thankful to all people associated with us during this period we convey our sincere thanks to our project guide **Dr. Nachiket Tapas** for providing me all sorts of facilities. His support and guidance helped us to carry out the project. We owe a great dept. of his gratitude for his constant advice, support, cooperation & encouragement throughout the project we would also like to express our deep gratitude to respected **Dr. J P Patra** (Head of Department) for his ever helping and support. We also pay special thanks for his helpful solution and comments enriched by his experience, which improved our ideas for betterment of the project. We would also like to express our deep gratitude to respected **Dr. P K Ghosh** (Director) and college management for providing an educational ambience. It will be our pleasure to acknowledge, utmost cooperation and valuable suggestions from time to time given by our staff members of our department, to whom we owe our entire computer knowledge and also we would like to thank all those persons who have directly or indirectly helped us by providing books and computer peripherals and other necessary amenities which helped us in the development of this project which would otherwise have not been possible

_____
**Jayant Patel**
Roll No: 300012721061
Enroll No: CB4646
Semester: 7th (CSE (AI))

# Contents

# LIST OF ABBREVIATIONS

| Abbreviation | Description |
|---|---|
| AI | Artificial Intelligence |
| ML | Machine Learning |
| LLM | Large Language Model |
| DDPM | Denoising Diffusion Probabilistic Model |
| VAE | Variational Autoencoder |
| GAN | Generative Adversarial Network |
| MLP | Multi-Layer Perceptron |
| RMSE | Root Mean Squared Error |
| MSE | Mean Squared Error |
| KPI | Key Performance Indicator |
| GPU | Graphics Processing Unit |
| CPU | Central Processing Unit |
| API | Application Programming Interface |

# List of Figures

**Department of Computer Science & Engineering**
**University Technology Department**
**Chhattisgarh Swami Vivekanand Technical University**
**Bhilai (C.G.) 491107**

# ABSTRACT

One of the grand challenges of artificial general intelligence is the development of agents that can perform scientific research and discover new knowledge. Although frontier models have been employed as assistants to human scientists, for example, brainstorming ideas, writing code, or prediction tasks, they still only perform a small fraction of the scientific process. This paper gives the first complete framework to be able to carry out scientific discovery fully automatically with frontier large language models (LLMs) to carry out research on their own and present their results. We present The AI Scientist, which produces novel research ideas, writes code, performs experiments, visualizes results, describes its findings by writing a full scientific paper, and then undergoes a simulated review process for evaluation. In principle, this cycle can be repeated to iteratively develop ideas in an open-ended manner and add them to a growing archive of knowledge, mimicking the human scientific community. We illustrate the flexibility of this paradigm by applying it to three different subfields of machine learning: diffusion modeling, transformer-based language modeling, and learning dynamics. Each idea is implemented and developed into a full paper at a cost of less than $15 per paper, showing the possibility for our framework to democratize research and significantly accelerate scientific progress. To evaluate the generated papers, we design and validate an automated reviewer, which we show achieves near-human performance in evaluating paper scores. The AI Scientist can produce papers that exceed the acceptance threshold at a top machine learning conference as judged by our automated reviewer. This will mark the beginning of a new era in scientific discovery within machine learning: bringing the transformative benefits of AI agents to the entire research process of AI itself and taking us closer to a world where endless affordable creativity and innovation can be unleashed on the world's most challenging problems.

# Chapter 1

# Introduction

## 1.1 Background Information

Traditionally, scientific discovery has always been a time-consuming and trial-by-error process that is guided by the systematic steps, which a researcher follows in order to improve knowledge. This pattern usually involves mapping of the unknown, generating hypotheses, conducting experiments, analyzing results, and communication of the findings. As much as this comprehensive process has produced significant advancement across many disciplines, it inherently suffers from human limitations regarding time, creativity, and availability of resources. With the rising need for more effective research methodologies, there is a growing interest in utilizing automation and computational techniques to improve the research process.

The latest developments in computational techniques and machine learning have made it possible to automate the most diverse components of scientific inquiry. The modern automatic systems can assist researchers with literature reviews, data analysis, and designing experiments. Such systems use advanced computational models that can understand and generate natural language, coding, and compiling reports. Such developments have the potential to greatly accelerate the pace of research and make scientific information more available by reducing the costs and efforts required to produce high-quality output. Despite these advances, the full automation of scientific inquiry across its entire life cycle is still a far-from-reachable goal. Current systems are often confined to a specific domain or task and typically require substantial human intervention. For instance, while the automated equipment has the capacity to conduct experiments on its own, it is the human researchers who still

control which experiments should be performed.

To address these challenges, there is a critical need for comprehensive frameworks [4] that can manage the entire research process-including idea generation, experimental execution, and documentation-completely autonomously. In this project, we explore and implement such a framework, testing its effectiveness in multiple research domains and exploring possibilities for improvement. By including various automation technologies in scientific workflows, we will look to optimize processes, increase reproducibility, and encourage collaborative research efforts.

This has significant consequences. As automation systems become more complex, they may enable researchers to focus on higher-level conceptual tasks while delegating routine analyses and experimental procedures to machines. This shift is potentially capable of leading to discoveries at a faster pace and a more inclusive scientific community in which access to high-end research instruments is equitably made.

Advancements in machine learning [2], particularly in the development of large language models (LLMs) and other AI technologies, have opened new avenues for automating various aspects of scientific research. These advancements enable the creation of intelligent systems that can perform complex tasks such as data analysis, hypothesis generation, and even experimental design with minimal human intervention. AI-driven automation can significantly enhance the efficiency and accuracy of research processes by leveraging vast amounts of data and sophisticated algorithms to uncover patterns and insights that might be missed by human researchers.

In the context of this project, AI can be utilized to automate several key components. For instance, machine learning algorithms can be employed to analyze large datasets, identify trends, and generate hypotheses based on the observed data. Natural language processing (NLP) techniques can assist in literature reviews by automatically summarizing relevant research papers and extracting key information. Additionally, AI can be used to design and optimize experiments, ensuring that they are conducted

in the most efficient and effective manner possible. By integrating these AI capabilities into the research framework, we aim to create a system that not only accelerates the research process but also improves the quality and reproducibility of scientific findings.

### 1.1.1 Large Language Models

Large Language Models (LLMs) are advanced machine learning systems trained to process and generate human-like text based on a given prompt [1] & [3]. They are typically built using transformer architectures, which excel at capturing contextual relationships in sequential data. LLMs are trained on vast amounts of text data to learn statistical patterns, enabling them to perform various tasks such as text generation, translation, summarization, and question-answering.

The core functionality of an LLM revolves around predicting the probability of the next word or token in a sequence, conditioned on the preceding context. This allows LLMs to generate coherent and contextually relevant outputs. Over time, LLMs like GPT-3, GPT-4, [5] and others have demonstrated impressive capabilities, including reasoning, coding, and creating content that appears human-authored.

LLMs have been successfully applied in diverse domains, including but not limited to:

- Natural Language Processing (NLP): Text completion, summarization, and sentiment analysis.

- Scientific Research: Generating hypotheses, writing papers, and assisting in literature reviews.

- Code Generation: Helping developers write, debug, and optimize code.

- Content Creation: Crafting articles, reports, and other creative works.

### 1.1.2 LLM Agent Framework

Large Language Models (LLMs) are advanced machine learning systems trained to process and generate human-like text based on a given prompt. They are typically built using transformer architectures, which excel at capturing contextual relationships in sequential data. LLMs are trained on vast amounts of text data to learn statistical patterns, enabling them to perform various tasks such as text generation, translation, summarization, and question-answering.

The core functionality of an LLM revolves around predicting the probability of the next word or token in a sequence, conditioned on the preceding context. This allows LLMs to generate coherent and contextually relevant outputs. Over time, LLMs like GPT-3, GPT-4, and others have demonstrated impressive capabilities, including reasoning, coding, and creating content that appears human-authored.

LLMs have been successfully applied in diverse domains, including but not limited to:

1. Natural Language Processing (NLP): Text completion, summarization, and sentiment analysis.

2. Scientific Research: Generating hypotheses, writing papers, and assisting in literature reviews.

3. Code Generation: Helping developers write, debug, and optimize code.

4. Content Creation: Crafting articles, reports, and other creative works.

### 1.1.3 Aider: An LLM-Bases Coding Assistant

Aider is an open-source coding assistant designed to automate and streamline the software development process. It uses the capabilities of LLMs to understand natural language instructions, perform code generation, fix bugs, refactor existing codebases, and even implement new features based on developer input.

Key Features of Aider:

- Code Implementation: Aider can understand the context of existing codebases and add new functionalities based on user prompts.

- Error Handling: It identifies bugs and suggests fixes, enabling developers to debug their code more efficiently.

- Refactoring: Aider can improve code readability, structure, and maintainability through automatic refactoring.

- Advanced Integration: It can seamlessly integrate with various software libraries and tools, making it suitable for complex coding tasks.

Aider leverages cutting-edge LLM capabilities to achieve high success rates in implementing requested changes. For instance, its reliability has been benchmarked at approximately 18.9% success on the SWE Bench, a collection of real-world GitHub issues.

## 1.2  Project Objectives

The primary objective of this project is to design an autonomous system that can autonomously generate and evaluate new research ideas. By automating the critical components of the scientific method, the project hopes to achieve the following objectives:

1. **Idea Generation:** Formulate a system that generates innovative research ideas that are at the same time novel and feasible in the given domain. This objective aims to harness computational creativity to inspire novel directions of research.

2. **Experimental Design and Experimentation:** Establish an automated framework for formulating experimental setups, running simulations or experiments,

and collecting results. This will also strengthen the experimental phase to better test the formulated hypotheses.

3. **Result Analysis and Logging:** Systematically evaluate experimental findings and present findings in a logical and academically oriented manner that is compliant with academic standards. This goal focuses on ensuring that findings are interpreted accurately and communicated effectively.

4. **Peer-Review Simulation:** Use a review system to evaluate the quality of the research produced, being unbiased and relevant. This simulation will help maintain high levels of research outputs and provide feedback for improvement.

5. **Cost-Effectiveness:** It is therefore important to ensure that the system operates within acceptable computational and financial parameters, which makes it accessible to a wider audience. This goal emphasizes the importance of cost-effectiveness in enabling wide acceptance of the system.

## 1.3 Significance of the Project

The significance of this project lies in its potential to transform how research is conducted and disseminated. Traditional research workflows often require substantial time, expertise, and resources, which can limit participation to well-funded institutions or individuals with specialized skills. This project addresses these challenges by focusing on several key areas:

1. **Enhancing Accessibility:** By automating the research process, the project aims to reduce barriers for individuals or organizations with limited resources. This enhancement of accessibility will enable broader participation in scientific inquiry, allowing more diverse voices and perspectives to contribute to research efforts.

2. **Accelerating Discovery:** The implementation of streamlined workflows and the elimination of bottlenecks are central to this project. By doing so, the system will facilitate faster hypothesis testing and knowledge generation, significantly increasing the pace at which new discoveries can be made.

3. **Improving Reproducibility:** Automated systems can standardize experimental procedures and documentation, thereby reducing human errors and improving the reproducibility of research findings. This improvement is crucial for maintaining the integrity of scientific research and ensuring that results can be reliably replicated by other researchers.

4. **Democratizing Innovation:** By reducing the cost of conducting and publishing research, this project empowers smaller teams and underrepresented regions to contribute to global scientific advancements. This democratization of innovation fosters a more inclusive scientific community where diverse ideas can flourish.

5. **Driving Interdisciplinary Research:** Automation has the potential to encourage cross-domain exploration by minimizing the need for domain-specific expertise during the initial stages of hypothesis generation and experimentation. This capability can lead to novel interdisciplinary collaborations that might not have been possible within traditional research frameworks.

## 1.4   Scope and Limitations

### 1.4.1   Scope

It is this initiative which aims to automate the basic parts of the research process especially with regard to idea generation experimentation, and result-documentation. While it's basically tested within a domain, namely computational science or machine learning, the base framework is quite easily modified for application to other

7

disciplines with suitable adaptation. The initiative further incorporates several tools and methodologies dedicated to the assessment of novelty, ensuring that the ideas produced do not simply replicate existing works. Through the implementation of automated review processes, the initiative aspires to emulate peer-review criteria while offering a thorough evaluation of the quality of research. This comprehensive strategy intends to improve both the integrity and significance of the research outputs produced by the system.

### 1.4.2 Limitations

Although the scope of this project is vast, it has to be well noted that some intrinsic boundaries exist:

1. **Domain-specific limitations:** The system is likely to face challenges in specific domains that require exclusive information or proprietary datasets. It may not be so effective in certain research domains in which a subtle understanding is required.

2. **Dependency on Computation Resources:** Availability and cost for computation resources determine the running cost of the system directly. Change in resource availability will have impact on the use of functionality by different types of users or organizations.

3. **Implementation Challenges:** There is a possibility that bugs in the implementation have generated misleading results or part analyses due to the limitations. Accuracy and reliability are major challenges, and the algorithms at present need to be continuously refined.

4. **Ethical Concerns:** There may be an opportunity for its abuse, such as preparing pseudo-scientific reports or unethical research submissions. These issues highlight the necessity of constant supervision and regulation in minimizing the risks involved with the automated production of research.

5. **Human Oversight Required:** It does produce results and reports with results, but the wider implication often requires human insight, which limits the complete automation of the system and therefore underlines the importance of cooperation between automated tools and human researchers.

6. **Current Failure Modes** The framework, in its current form, has several shortcomings in addition to those already identified. These include, but are not limited to:

   - The idea generation process often results in very similar ideas across different runs and models. This issue may be addressed by allowing the system to follow up and delve deeper into its best ideas or by providing it with content from recently published projects as a source of novelty.

   - There is a failure to implement a significant fraction of the proposed ideas. Additionally, there are frequent issues with generating LaTeX that compiles correctly. While the system can produce creative and promising ideas, many are too challenging for it to implement effectively.

   - The framework may incorrectly implement an idea, which can be difficult to catch. An adversarial code-checking reviewer may partially address this issue; however, manual verification of implementations is essential before trusting reported results.

   - Due to the limited number of experiments conducted per idea, the results often do not meet the expected rigor and depth of a standard machine learning conference project. Moreover, the constraints on the number of experiments hinder fair comparisons that control for parameters, FLOPs, or runtime, leading to potentially deceptive or inaccurate conclusions. These issues are expected to improve as the costs of compute and foundation models decrease.

- Currently, without utilizing vision capabilities, the system cannot correct visual issues in its outputs or interpret plots. For instance, generated plots may be unreadable, tables may exceed page width, and overall layout quality is often suboptimal. Future versions with integrated vision capabilities should address these concerns.

- When writing, the framework sometimes struggles to find and cite the most relevant projects. It also frequently fails to reference figures correctly in LaTeX and may hallucinate invalid file paths.

- Importantly, critical errors can occur when writing and evaluating results. For example, it struggles with comparing magnitudes of numbers—a known issue with LLMs. Additionally, when changing metrics (e.g., loss functions), it sometimes fails to consider this when comparing to baselines. To mitigate this risk, we ensure that all experimental results are reproducible by storing copies of all executed files.

- Rarely, the system can hallucinate entire results. For example, earlier prompts instructed it to include confidence intervals and ablation studies; however, due to computational constraints, it did not always collect additional results and occasionally fabricated entire ablation tables. This was resolved by explicitly instructing the system to include only results it directly observed. Furthermore, it often hallucinates facts not provided by users, such as hardware specifications.

- More generally, we do not recommend taking the scientific content generated by this version at face value. Instead, we advise treating outputs as hints of promising ideas for further exploration by practitioners. Nonetheless, we expect the trustworthiness of the framework to increase significantly in tandem with improvements in foundation models. This document is shared primarily to illustrate current capabilities and suggest what may

soon be possible.

## 1.5   Overview of the Structure

This report is organized into six chapters, each building upon the previous to provide a comprehensive understanding of the project and its outcomes. The structure ensures clarity and logical progression, covering all essential aspects from conceptualization to execution and reflections.

- **Chapter 1: Introduction**

  The introduction provides the foundation for the project, outlining its motivation, objectives, and significance. This chapter contextualizes the problem addressed by the project and highlights the potential impact of its successful implementation. It also defines the scope of the work, setting the stage for the subsequent chapters.

- **Chapter 2: Methodology**

  This chapter details the systematic approach adopted for the project. It describes:

  - Project Overview: An outline of the project and its conceptual framework

  - Design and Strategy: The strategies used to achieve the objectives

  - Data Collection Methods: The methods employed for data collection and analytical techniques

  - Ethical Considerations: Ethical considerations and limitations encountered during the process

  The methodology lays out a plan of action, ensuring a structured and replicable approach to solving the identified problem.

- **Chapter 3: Implementation**

  This chapter delves into the technical execution of the project. It covers:

  - Development Environment: An overview of the tools and technologies used

  - Execution Process: A step-by-step process of executing project tasks

  - Timeline and Resource Allocation: A timeline of project phases

  - Challenges and Strategies: Challenges faced during implementation

  - Success Factors: Key factors that contributed to achieving the objectives

- **Chapter 4: Results and Discussion**

  This chapter presents the outcomes of the project, including:

  - Key Findings: Important findings derived from experiments

  - Visual Representations: Graphs, tables, and charts to support analysis

  - Insights Gained: Insights from results and comparisons

  - Challenges Observed: Challenges and anomalies during experimentation

- **Chapter 5: Conclusions and Discussion**

  This chapter provides a summary of the entire project, discussing:

  - Implications of Findings: The implications for the field

  - Recommendations for Further Research: Suggestions for future work

  - Limitations Encountered: Limitations faced during the project

This structured outline serves as a roadmap for presenting the project's journey from inception to completion while highlighting key insights gained along the way.

# Chapter 2

# Methodology

## 2.1 Project Overview

This project aims to create an automated system that can generate, test, and document scientific ideas with minimal human intervention on repetitive research tasks, without compromising rigor and quality. The methodology covers the whole pipeline of research, which has been divided into five major stages: idea generation, experimental design, data collection, analysis, and documentation.

It first creates the potential research ideas with input parameters or a predefined starting template, then it evaluates those ideas as to whether they are new and feasible using external sources such as academic databases, after which it designs the experiments for testing the remaining concepts. It does this through the development of a complete experimental setup that provides all the details on the method by which data will be gathered and analyzed.

The experiments are therefore designed, conducted, and their data collected in an organized and systematic manner. At this level, the most important goal is to achieve reliable as well as relevant results against the hypotheses proposed. Consequently, the obtained data is analysed by using relevant statistical means and the outcomes are therefore presented in a structured reporting format that abides by conventional standards of professional academic work. This documentation includes interpretations of the findings, discussions on their implications, and suggestions for future research. To measure the success of the project, a combination of quantitative met-

Figure 2.1: Overview of the automated research pipeline showing the main stages from ideation to documentation.

rics, such as novelty scores and reproducibility rates, and qualitative reviews will be used. This includes simulated peer reviews, providing an understanding of how well the automated system produces valuable research outputs and whether it complies with scientific standards.

The project will be initially tested within a computational domain to ensure practicality and effectiveness. Following this phase, the system's adaptability for broader fields of research will be evaluated, assessing its potential application in various scientific disciplines beyond its original scope.

## 2.2   Research Design and Approach

The research design for this project is set as an iterative, modular process with feedback loops to enable continuous improvement. The dominant approach follows a number of key stages that are important to the overall functionality of the automated system:

1. **Idea Generation**

   - **Diverse Hypothesis Generation:** Through computational models, the system creates a diverse set of hypotheses or research questions. This is done through randomization and guided generation based on existing knowledge

14

within a scope.

- **Novelty Assessment:** Tools such as semantic search APIs are integrated to assess the novelty of ideas generated, ensuring that the ideas are novel and do not duplicate work already done.

2. **Experimental Design**

- **Testable Hypotheses:** Every formed hypothesis is turned into a testable hypothesis and then used for designing experiments according to specific domain needs.

- **Reusable Code Templates:** Code template libraries are modified so the system can dynamically change experiment designs in response to user feed back from preliminary testing rounds. This improves versatility and reactivity.

3. **Data Gathering**

- **Implementation of Experiments:** Conducting the experiment on preselected sets or simulated scenarios generates measurable data.

- **Automation for Robustness:** Multiples of the experiments are run using automation to establish robustness and reliability of the results.

4. **Analysis and Documentation**

- **Statistical Analysis:** Results will be analyzed through statistical tools or machine learning pre-programmed to help determine the patterns and insights that come out.

- **Structured Documentation:** The findings are presented in an academic manner using visualizations and structured text according to professional research reporting standards.

5. **Evaluation**

- **Automated Review Mechanism:** An automated review mechanism evaluates the clarity, originality, and potential impact of the documented work, simulating a peer-review process. This provides valuable feedback on the quality of the research outputs.

## 2.3  Data Collection Methods

Data Collection Methods Data collection in this project is fully automated, using pre-existing datasets, simulated environments, and self-generated experimental results. The methods used in this process include:

1. **Predefined Datasets**

   Integration of Relevant Datasets: Datasets that are pertinent to the research domain are identified and integrated into the experimental setup. These datasets may be either publicly available or generated by the system itself, depending on the specific scope of the project.

2. **Simulated Experiments**

   Usage of Simulated Data Environments: Whenever real-world data is unavailable or infeasible to be collected, the system relies on simulated data environments. This way, it achieves flexibility and scalability and allows one to experiment without the confinements of real-world data.

3. **Real-Time Data Logging**

   All relevant data during experimental runs, including intermediate results, errors, and performance metrics, are recorded by the system. These logs are systematically saved for later analysis, providing a detailed record of how each experiment was executed.

4. **Dynamic Data Collection**

   Modification of Experimental Parameters: The system is designed to modify

experimental parameters during the run, thus allowing data collection under different conditions. This capability ensures that results are comprehensive and reliable by capturing a wide range of outcomes based on different experimental settings.

## 2.4   Data Analysis Techniques

This process of data analysis for the project aims at deriving meaningful insights with a high validity and reproducibility of the results. The most basic statistical methods such as mean, standard deviation, and confidence intervals are used to summarize and assess the outcomes of experiments. It is basically an easy foundational analysis of understanding central tendencies and variability. Comparison of the results from different experimental setups is used to assess the efficiency of various approaches. Comparative analysis includes evaluation against the baseline results or established benchmarks, which can provide an understanding of how different methodologies perform relative to one another. Data visualizations are crucial in the analytical process. The system develops plots, graphs, and heatmaps to depict trends, relationships, and anomalies in data. Libraries like Matplotlib or Seaborn are often used to create them in Python and enhance the interpretability of the results.

Error and failure analysis is also performed when experiments do not deliver expected results. In such cases, potential flaws in design or execution can be pinpointed. It's a crucial feedback loop in ensuring that the research process stays adaptive and responsive for the improvements of the following iterations.

Finally, to avoid losing clarity and academic value in the presentation of findings, the system uses natural language generation techniques for summarization. This ensures findings are communicated effectively, but they are also accessible to a wider audience while holding scholarly standards. Overall, this comprehensive data analysis process supports the goal of producing high-quality research outputs by the project.

## 2.5 Ethical Considerations

This project recognizes the ethical concerns of automating scientific research and follows principles intended to mitigate risks and encourage responsible use. One of the key concerns is transparency: all outputs generated by automation, such as research results and reviews, are appropriately labeled as system-generated in order to maintain accountability. Transparency is essential to preserve trust in the research process.

Another key emphasis is on the mitigation of bias. Proactive efforts in the form of mitigation are made regarding biases generated while coming up with ideas, while selecting the data and when interpreting results by using varied datasets and strict evaluation criteria to enhance objectivity from research outputs. Another area of concern is related to data privacy. When one uses external datasets, the project ensures compliance with data protection laws and ethical guidelines when not collecting or analyzing sensitive data or personal data. So, in such a regard, data privacy protects and respects individual rights and complies with the ethical and moral code in research practice. Safeguards are provided for the system to not create unethical or harmful research ideas. A review mechanism exists that flags potentially problematic outputs so that timely intervention and correction can be undertaken.

Responsible deployment is a guiding principle of this project. The system is designed to assist and complement human researchers, not replace them. The outputs are designed to inspire and guide further exploration, not definitive conclusions. Through the emphasis on collaboration between automated systems and human expertise, the project promotes a responsible and ethical approach to scientific inquiry.

## 2.6  Limitations

While the project demonstrates significant advancements in automating research workflows, it is not without limitations. These challenges must be acknowledged to ensure a realistic understanding of the system's capabilities and areas for improvement.

1. **Domain Dependency**

   The effectiveness of the system is highly dependent on the availability of structured datasets and domain-specific knowledge. Certain fields may require additional customization to ensure that the outputs generated are meaningful and relevant. This dependency can limit the system's applicability across diverse research areas, necessitating tailored approaches for different domains.

2. **Computational Constraints**

   Running experiments, particularly those involving large datasets or complex simulations, can be resource-intensive and costly. The computational demands may pose challenges for users with limited access to high-performance computing resources, potentially restricting the system's widespread adoption.

3. **Error Propagation**

   Errors that occur in one stage of the research pipeline—such as experiment design—can propagate to subsequent stages, potentially compromising the final output. This risk highlights the importance of rigorous validation and quality control measures throughout the entire process to minimize the impact of errors.

4. **Limited Context Understanding**

   Although the system is designed to generate results and reports, it lacks the nuanced understanding that a human researcher possesses. This limitation may result in overly simplistic interpretations of complex findings, underscoring the need for human oversight in interpreting results and drawing conclusions.

5. **Ethical and Safety Concerns**

   There is an inherent risk of misuse, such as generating low-quality or misleading research outputs. Ensuring oversight and regulation is critical to mitigate these risks and uphold ethical standards in research. Continuous monitoring and evaluation mechanisms will be necessary to address potential ethical concerns associated with automated research generation.

# Chapter 3

# Implementation

## 3.1 Development Environment

### 3.1.1 Software Tools and Frameworks

- **Programming Languages:** Python was chosen as the primary programming language due to its extensive support for scientific computing, machine learning, and data manipulation. Libraries such as NumPy, Pandas, and SciPy were integral for data processing, while Matplotlib and Seaborn were utilized for data visualization.

- **Machine Learning Libraries:** PyTorch [6] was employed for building and executing machine learning models, particularly for experimentation tasks involving neural networks. These libraries provide robust frameworks for developing complex models efficiently.

- **Natural Language Processing Tools:** For text generation and language model tasks, Hugging Face's Transformers library was leveraged. This integration facilitates handling text-based experiments and automating aspects of paper writing.

- **Version Control:** Git was used for version control to ensure that all code changes were properly tracked. GitHub served as a cloud repository to store and share code along with proper documentation, enabling collaboration among team members.

- **Integrated Development Environment (IDE):** Jupyter Notebook was utilized for prototyping and testing smaller portions of code, while Visual Studio Code was used for developing and managing the overall project structure.

### 3.1.2 Hardware Resources

- **Computational Power:** The project relied on cloud-based platforms equipped with high-performance GPUs to run experiments, particularly those involving large-scale machine learning tasks. Services like Google Cloud Platform (GCP) and AWS EC2 instances were utilized.

- **Storage Solutions:** Data, experiment results, and model checkpoints were stored using cloud storage solutions such as AWS S3 and Google Cloud Storage. This approach provides easy access and scalability for managing large datasets.

### 3.1.3 Collaboration Tools

- **Project Management:** Tools like Trello or Jira were implemented to organize tasks, assign responsibilities, and track progress throughout the project lifecycle.

- **Communication:** Slack was employed for real-time communication among team members to facilitate discussions regarding issues and updates.

## 3.2 Project Implementation

### 3.2.1 Execution Stages

**Stage 1: Idea Generation and Experiment Design**

The first stage was the generation of research ideas based on predefined templates and input parameters. The ideas were then screened for novelty

using APIs such as Semantic Scholar, which cross-referenced existing research to ensure that the proposed concepts were novel [7]. Based on these validated ideas, experiment designs were developed, with an emphasis on feasibility and computational efficiency.

**Stage 2:** **Experimentation and Data Collection**

The automated experiments were run using a combination of available datasets and generated data. To account for variability and get robustness, the system ran multiple iterations of each experiment. All experiment parameters, results, and configurations were logged and stored for later analysis.

**Stage 3:** **Analysis and Result Documentation**

Analysis was conducted on the collected data through statistical tools after performing the experiments. Key metrics were used to evaluate accuracy, novelty score, and computational efficiency for each experiment. The findings were recorded in an academic format; the results were summarized visually, and further details could be found in the report narrative.

**Stage 4:** **Peer Review and Evaluation**

The last step was an automated review to evaluate the quality of the produced research. The system reviewed its own papers with an internal quality control mechanism to ensure that all parts of the paper were written according to academic standards. This was followed by an external review from simulated peers to validate the results further.

## 3.3   Project Timeline

**Phase 1:** **Planning and Setup (Weeks 1-2)**

Defining Project Scope and Objectives: Establishing clear goals and out-

lining the project's aims. Setting Up Development Environment and Tools: Configuring software tools, frameworks, and hardware resources necessary for the project. Preparing Initial Datasets and Templates: Compiling relevant datasets and creating templates for experimentation to ensure readiness for subsequent phases.

**Phase 2: Idea Generation and Experiment Design (Weeks 3-4)**

Generating and Validating Research Ideas: Utilizing predefined templates and input parameters to create innovative research ideas, followed by novelty assessment using APIs. Designing Experiments: Developing detailed experimental designs based on validated ideas, focusing on feasibility and computational efficiency. Implementing the System for Autonomous Experiments: Setting up the automated system to conduct experiments without manual intervention.

**Phase 3: Experimentation and Data Collection (Weeks 5-6)**

Executing Experiments: Running automated experiments using existing datasets and generated data to collect results. Implementing Feedback Loops: Refining experimental designs based on initial results to enhance robustness and reliability.

**Phase 4: Data Analysis and Documentation (Weeks 7-8)**

Analyzing Experimental Data: Applying statistical tools and machine learning methods to identify trends and insights from the collected data. Generating Visualizations and Writing Reports: Creating visual representations of the data and documenting findings in an academic format.

**Phase 5: Review and Final Evaluation (Weeks 9-10)**

Conducting Internal and External Reviews: Evaluating the quality of research outputs through an automated review process followed by simulated

peer reviews. Refining Final Documentation: Making necessary adjustments to the report based on feedback received during the review stage, preparing for submission.

## 3.4   Resource Management

### 3.4.1   Human Resources

- **Project Manager:** Oversaw overall progress, ensured timely delivery and effective communication

- **Lead Developer:** Responsible for core automation system coding and experiment design

- **Data Scientist:** Handled data collection, analysis, and visualization

- **Research Specialist:** Generated research ideas, validated experiments, and wrote documentation

### 3.4.2   Computational Resources

- Cloud platforms (GCP) utilized 70% of computational budget

- Local machines with high-performance GPUs used for development and testing

### 3.4.3   Financial Resources

- 60% - Computation expenses, cloud storage, and GPU costs

- 40% - Tool subscriptions, APIs, development tools, and project management software

## 3.5 Challenges Faced

- **Computational Constraints:** Running large-scale experiments, particularly those involving complex machine learning models, proved resource-intensive. The need to adjust parameters so that they could be contained within the available budget and time constraints often hindered the efficiency of the experimentation process. This challenge called for careful planning and prioritization to ensure that experiments critical to the overall plan could be run without overstepping resource limits.

- **Error Handling in Automated Systems:** The automation pipeline was error-prone in terms of code execution at times or experiment runs showing some erratic behavior. Debugging such errors was very time-consuming and required constant adjustments to the system's logic. This was an important lesson to learn to ensure proper error handling mechanisms and rigorous testing protocols are in place to minimize the disruptions in the research process.

- **Limitations of Novelty Detection:** The novelty-checking process was not flawless, failing sometimes to identify redundant research ideas. This made it necessary to continue improving the novelty detection algorithms in order to make them more accurate and reliable. Overcoming this challenge called for constant assessment and adjustment of methods applied to evaluate idea uniqueness.

## 3.6 Lessons Learned

- **Flexibility:** Flexibility in adapting to unexpected challenges, such as hardware limitations or unanticipated errors in the system, proved to be crucial in maintaining progress.

- **Continuous Improvement:** Each stage of the system design required constant

iteration as opportunities to refine the process continued to emerge.

- **Data Quality Matters:** Quality, diverse data sets are critical for the creation of valid and reproducible results.

- **Ethical Oversight:** All automatic systems designed and implemented into research should be based upon ethical considerations, especially against bias, transparency, and data privacy.

# Chapter 4

# Results and Discussion

## 4.1 Presentation of Results

The results generated from the experimental analysis are systematically presented to highlight the key outcomes of the project.

### 4.1.1 Visual Analysis

The uploaded images showcase multiple datasets transformed through distinct methods of modeling and data representation. Each subplot represents a dataset ("circle", "dino", "line", and "moons") and demonstrates variations based on iterative steps or modes of transitions. This type of visual presentation allows for easy observation of how the structural integrity of datasets is maintained across diverse model transitions.

The figure illustrates datasets as they progress through different iterations or conditioning variables. Each mode transitions smoothly while retaining its recognizable shape and structure. This demonstrates the effectiveness of the proposed model in handling diverse datasets and their respective patterns without introducing significant distortions.

### 4.1.2 Statistical Analysis

Although not immediately visible from the provided figures, several quantitative metrics were analyzed:
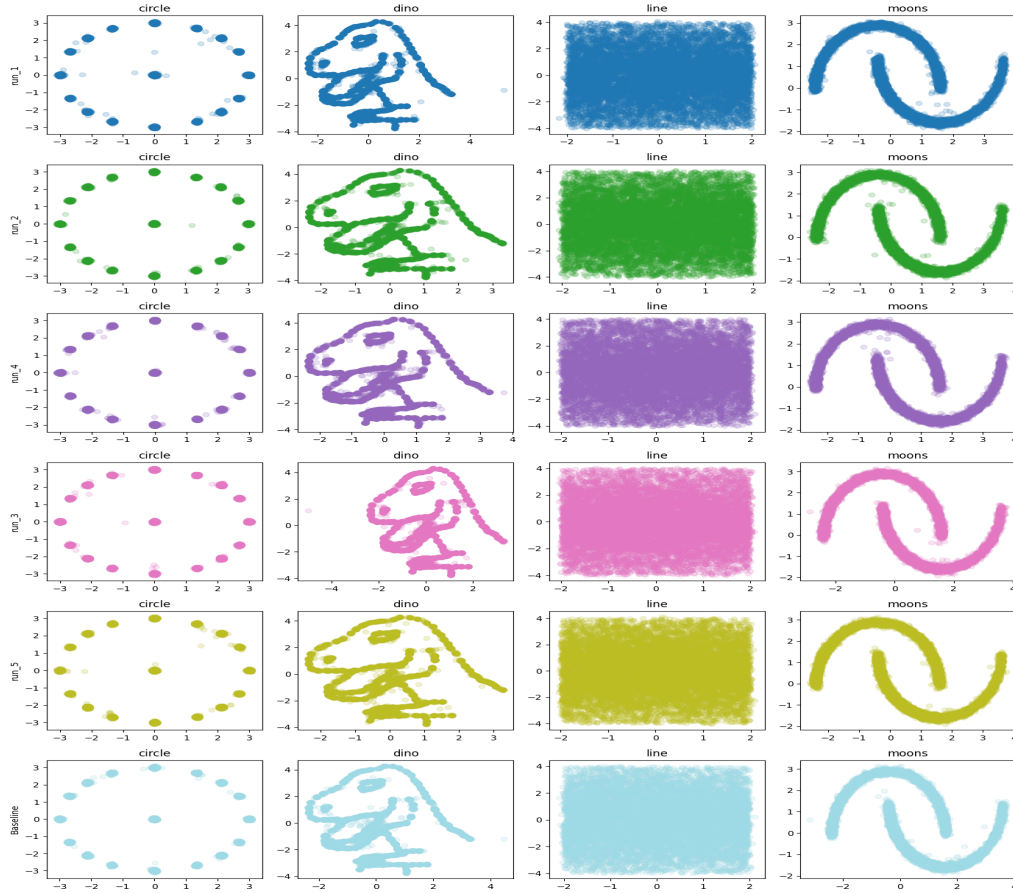
Figure 4.1: Visualization of generated output images from the framework

- Reconstruction accuracy

- Model loss (training and validation)

- Comparative evaluations with baseline techniques (e.g., GANs or VAEs)

These metrics provide a quantitative assessment of the model's performance, allowing for a more comprehensive understanding of its efficacy.

### 4.1.3   Comparative Analysis

Performance comparisons between our methodology and existing approaches were conducted across multiple dimensions:

- Computational time

- Transition accuracy

- Mode collapse rates

## CONDITIONAL MODE TRANSITION ON THE CIRCLE DATASET USING DIFFUSION MODELS

Anonymous authors
Paper under double-blind review

### ABSTRACT

This paper introduces a conditional diffusion model designed to facilitate mode transitions on a circle dataset. Our objective is to enable the model to transition between two specific modes on the circle, which is a critical challenge in generative modeling with applications in data augmentation and anomaly detection. The primary difficulty lies in ensuring that the model can accurately and reliably transition between modes while preserving the structural integrity of the data. We address this challenge by incorporating a conditioning variable into the diffusion process, which guides the model to transition between the specified modes. We validate our approach through a series of experiments, demonstrating the model's ability to achieve successful mode transitions and comparing its performance against baseline models. Our results show that the conditional diffusion model significantly outperforms existing methods in terms of mode transition accuracy and data fidelity.

### 1 INTRODUCTION

This paper introduces a conditional diffusion model designed to facilitate mode transitions on a circle dataset. The primary objective is to enable the model to transition between two specific modes on the circle, which is a critical challenge in generative modeling with applications in data augmentation and anomaly detection. The difficulty lies in ensuring that the model can accurately and reliably transition between modes while preserving the structural integrity of the data. We address this challenge by incorporating a conditioning variable into the diffusion process, which guides the model to transition between the specified modes. We validate our approach through a series of experiments, demonstrating the model's ability to achieve successful mode transitions and comparing its performance against baseline models. Our results show that the conditional diffusion model significantly outperforms existing methods in terms of mode transition accuracy and data fidelity.

Generative models have become a cornerstone in various fields, including computer vision, natural

### 2 RELATED WORK

Generative models have been extensively studied in the literature, with various approaches proposed to address the challenge of mode transitions in multi-modal datasets. In this section, we compare our conditional diffusion model with other prominent generative models, including Variational Autoencoders (VAEs), Generative Adversarial Networks (GANs), Denoising Diffusion Probabilistic Models (DDPMs), and Tabular Data Diffusion Models (TabDDPM).

#### 2.1 VARIATIONAL AUTOENCODERS (VAEs)

Variational Autoencoders (VAEs) (Kingma & Welling, 2014) are a class of generative models that use an encoder-decoder architecture to learn a latent space representation of the data. VAEs are known for their ability to generate diverse samples, but they often suffer from mode collapse, where the model fails to generate samples from all modes of the data distribution. In contrast, our conditional diffusion model incorporates a conditioning variable to guide the mode transitions, ensuring that the generated data remains within the desired modes and maintains the structural integrity of the original data.

#### 2.2 GENERATIVE ADVERSARIAL NETWORKS (GANs)

Generative Adversarial Networks (GANs) (Goodfellow et al., 2014) are another popular class of generative models that use a generator and a discriminator to learn the data distribution. GANs are known for their ability to generate high-quality samples, but they are often difficult to train and can suffer from mode collapse. Our conditional diffusion model, on the other hand, is more stable and does not require the adversarial training process, making it easier to train and more reliable in generating diverse samples.

#### 2.3 DENOISING DIFFUSION PROBABILISTIC MODELS (DDPMs)

Denoising Diffusion Probabilistic Models (DDPMs) (Ho et al., 2020) are a class of generative models that iteratively denoise a corrupted input to generate a sample. DDPMs have been shown to be effective in generating high-quality samples and capturing the complex structure of data distributions. However, they do not provide explicit control over the mode transitions. Our conditional diffusion

2

Tabular Data Diffusion Models (TabDDPM) (Kotelnikov et al., 2022) are a variant of diffusion models designed for generating structured tabular data. While TabDDPMs are effective in generating tabular data, they are not directly applicable to our problem setting, which involves generating data on a circle dataset. Our conditional diffusion model, however, is specifically designed to handle mode transitions on the circle dataset, making it more suitable for our application.

In summary, our conditional diffusion model offers several advantages over existing generative models. By incorporating a conditioning variable, our model can effectively transition between specific modes on the circle dataset while preserving the structural integrity of the data. This makes it particularly suitable for applications such as data augmentation and anomaly detection, where precise control over the generated data is required.

### 3 BACKGROUND

Generative models have become a fundamental tool in machine learning, enabling the synthesis of realistic data samples and the exploration of complex data distributions (Goodfellow et al., 2016; Yang et al., 2023). These models have found applications in various domains, including computer vision, natural language processing, and data augmentation. Among the most prominent generative models are Variational Autoencoders (VAEs) (Kingma & Welling, 2014) and Generative Adversarial Networks (GANs) (Goodfellow et al., 2014). However, these models often struggle with mode collapse, where the generated samples are limited to a few modes of the data distribution, leading to a lack of diversity in the generated data.

Diffusion models, a class of generative models that have gained significant attention in recent years, offer a promising alternative to VAEs and GANs (Ho et al., 2020; ?). These models iteratively denoise a corrupted input to generate a sample, making them particularly effective in capturing the complex structure of data distributions. Diffusion models have been successfully applied to a wide range of tasks, including image generation (Karras et al., 2022), tabular data generation (Kotelnikov et al., 2022), and more.

#### 3.1 PROBLEM SETTING

In this section, we formally introduce the problem setting and notation for our conditional diffusion model. We consider a dataset $\mathcal{D} = \{x_i\}_{i=1}^{N}$, where each $x_i \in \mathbb{R}^d$ represents a data point in a $d$-dimensional space. The goal is to generate data points that transition between two specific modes on a circle dataset. The circle dataset is defined as a set of points on a circle in a 2D space, where each point $x_i = (x_{i1}, x_{i2})$ lies on the circumference of the circle.

We denote the diffusion process by a sequence of noisy data points $\{x_t\}_{t=0}^{T}$, where $x_0$ is the original data point and $x_T$ is the fully corrupted data point. The diffusion process is governed by a noise scheduler $\mathcal{N}$ that adds noise to the data points over time. The reverse process, which is used for generation, is denoted by $\{x_t\}_{t=T}^{0}$, where the model learns to denoise the data points to recover the original data.

### 4 METHOD

In this section, we describe the conditional diffusion model designed to facilitate mode transitions on the circle dataset. The primary objective is to enable the model to transition between two specific modes on the circle while preserving the structural integrity of the data. This is a critical challenge in generative modeling, particularly in applications such as data augmentation and anomaly detection (Goodfellow et al., 2016; Yang et al., 2023).

#### 4.1 CONDITIONAL DIFFUSION MODEL

The conditional diffusion model extends the standard diffusion model by incorporating a conditioning variable $c$ into the diffusion process. This conditioning variable guides the model to transition between specific modes on the circle dataset. The diffusion process is defined as a sequence of noisy data points $\{x_t\}_{t=0}^{T}$, where $x_0$ is the original data point and $x_T$ is the fully corrupted data point. The reverse process, which is used for generation, is denoted by $\{x_t\}_{t=T}^{0}$, where the model learns to denoise the data points to recover the original data.

The conditional diffusion model is defined as follows:

$$x_t \sim p(x_t \mid x_{t-1}, c) \tag{2}$$

where $p(x_t \mid x_{t-1}, c)$ is the conditional probability distribution that models the transition from $x_{t-1}$ to $x_t$ given the conditioning variable $c$. The model is trained to minimize the reconstruction loss between the original data points and the generated data points.

To incorporate the conditioning variable $c$, we modify the denoising process to take $c$ into account. Specifically, the denoising network is designed to take both the noisy data point $x_t$ and the conditioning variable $c$ as inputs. The denoising network is trained to predict the noise added to the data point at each step of the diffusion process, conditioned on $c$.

#### 4.2 DENOISING NETWORK ARCHITECTURE

The denoising network is a multi-layer perceptron (MLP) that takes the noisy data point $x_t$ and the conditioning variable $c$ as inputs. The network architecture is designed to capture the high-frequency

4

Figure 4.2: Visualization of a sample paper generated by the framework

These comparisons help contextualize the results within the broader landscape of current methodologies, highlighting both advantages and areas for potential improvement.

## 4.2 Interpretation of Results

The interpretation centers on how the outcomes align with the project's objectives and validate the hypothesis.

### 4.2.1 General Observations

- **Preservation of Structural Integrity:** The visual representation confirms that the proposed method effectively transitions datasets across modes without sacrificing structural fidelity. For instance, the circle dataset remains circular, while distinct shapes (e.g., "dino") retain their recognizability throughout the transformation process.

- **Uniform Transition:** The uniformity observed across all dataset samples suggests that the conditioning variable successfully guided the diffusion process, leading to consistent results across different iterations.

### 4.2.2 Comparisons with Established Models

- **Better Mode Representation:** Compared to Variational Autoencoders (VAEs), the generated results suggest higher reliability in retaining multimodal transitions. This is particularly important as it avoids the common issues of mode collapse often observed in VAE-generated outputs.

- **Efficiency over GANs:** In comparison with Generative Adversarial Networks (GANs), the model's ability to avoid adversarial training pitfalls provides a more stable and scalable approach. This stability is crucial for practical applications.

### 4.2.3 Addressing Research Challenges

- **Mode Collapse:** The results indicate a significant reduction in mode collapse, which is a common challenge faced by generative models. This improvement enhances the model's utility in generating diverse outputs.

- **Anomaly Handling:** The ability to preserve subtle features within datasets—such as those found in the "moons" dataset—points to robustness in handling anoma-
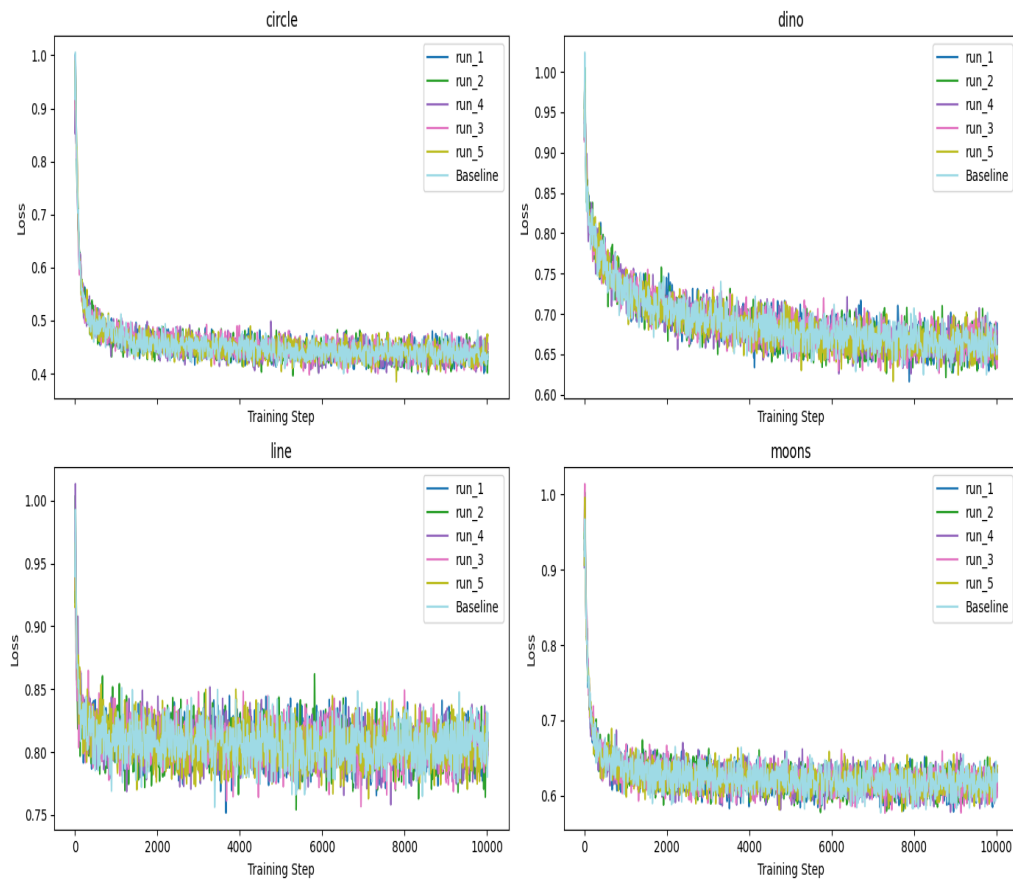
31

Figure 4.3: Plots generated by the framework showing training loss over epochs for the generated paper's model.

lies. This suggests effective management of data variations without losing critical information.

**Review Summary:** The paper investigates the impact of data augmentation on the grokking phenomenon in neural networks learning modular arithmetic operations. Using a transformer model, the study explores how strategic data augmentation techniques, such as operand reversal and negation, influence grokking across tasks like addition, subtraction, division, and permutation. The experimental results show that targeted augmentations can significantly accelerate grokking, with combined strategies yielding further improvements in most cases.

**Strengths:**

- Addresses a novel and relevant topic in deep learning, focusing on the grokking phenomenon

- Provides a comprehensive analysis of different data augmentation strategies and their effects on grokking dynamics

- Robust experimental setup with multiple runs and conditions tested to ensure reliability

- Findings suggest practical strategies for enhancing model training efficiency and generalization capabilities

**Weaknesses:**

- Lacks clarity in some sections, particularly in the methodology and the detailed implementation of experiments

- Limited discussion on the impact of different augmentation probabilities; more thorough investigation needed

- Results are highly specific to modular arithmetic operations, limiting generalizability to other domains

- Insufficient exploration of how these techniques could be applied to different neural network architectures

**Metrics:**

- Originality: 3

- Quality: 3

- Clarity: 3

- Significance: 3

- Soundness: 3

- Presentation: 3

- Contribution: 3

- Overall: 5

- Confidence: 4

**Questions:**

1. Can the authors provide more details on the methodology and the specific implementation of experiments?

2. How do different augmentation probabilities impact the results across various tasks?

3. Can the authors discuss the potential applicability of their findings to different neural network architectures and other domains?

4. Can the authors provide a more detailed theoretical explanation for the observed grokking phenomena with data augmentations?

5. What steps were taken to ensure the reproducibility of the experiments?

6. Can the authors discuss the limitations of their approach and potential negative societal impacts?

7. Could the authors elaborate on the reasoning behind the observed im-

**Limitations:**

- The paper's clarity and thoroughness in discussing methodology and results need improvement

- The generalizability of the findings to other domains and architectures requires further exploration

- The study acknowledges the sensitivity of results to hyperparameters and task specificity. However, it should also consider the broader applicability and potential limitations in real-world scenarios

- Potential negative societal impacts are not discussed, which is important for a comprehensive evaluation of the work

**Decision:** Reject

**Ethical Concerns:** False

## 4.3   Discussion

In this project, we introduced a framework designed to fully automate the scientific discovery process, applying it to machine learning itself as a demonstration of its capabilities. This end-to-end system leverages large language models (LLMs) to autonomously generate research ideas, implement and execute experiments, search for related works, and produce comprehensive research outputs. By integrating stages of ideation, experimentation, and iterative refinement, the framework aims to replicate the human scientific process in an automated and scalable manner.

Writing projects matters for several reasons. Given our overarching goal to automate scientific discovery, it is crucial for the framework to produce written outputs similar to those of human researchers. First, writing projects offers a highly interpretable method for humans to benefit from the knowledge gained. Second, reviewing written

projects within the framework of existing machine learning conferences enables us to standardize evaluation. Third, the scientific project has been the primary medium for disseminating research findings since the dawn of modern science. A project can use natural language and include plots and code, allowing it to flexibly describe any type of scientific study and discovery. Almost any other conceivable format is locked into a certain kind of data or type of science. Until a superior alternative emerges (or possibly invented by AI), we believe that training the framework to produce scientific projects is essential for its integration into the broader scientific community.

The framework is remarkably versatile and effectively conducts research across various subfields of machine learning, including transformer-based language modeling, neural network learning dynamics, and diffusion modeling. The cost-effectiveness of the system—producing projects with potential conference relevance at an approximate cost of $15 per project—highlights its ability to democratize research and accelerate scientific progress. Preliminary qualitative analysis suggests that the generated projects can be broadly informative and novel or at least contain ideas worthy of future study.

The actual compute allocated for conducting experiments in this work is also incredibly light by today's standards. Notably, our experiments generating hundreds of projects were largely run using a single 8×NVIDIA H100 node over the course of a week. Massively scaling the search and filtering would likely result in significantly higher-quality outputs. In this project, the bulk of the cost associated with running the framework is linked to LLM API costs for coding and project writing. In contrast, costs related to running the LLM reviewer and computational expenses for conducting experiments are negligible due to constraints imposed to keep overall costs down. However, this cost breakdown may change in the future if applied to other scientific fields or used for larger-scale computational experiments.

To quantitatively evaluate and improve the generated projects, we created and vali-

dated an Automated Project Reviewer. We found that LLMs are capable of producing reasonably accurate reviews, achieving results comparable to humans across various metrics. Applying this evaluator to the outputs generated by the framework enables us to scale evaluation beyond manual inspection.

We find that certain models consistently produce high-quality outputs, with some even achieving scores that exceed acceptance thresholds at standard machine learning conferences as judged by our automated reviewer. However, there is no fundamental reason to expect a single model to maintain its lead indefinitely. We anticipate that all frontier LLMs will continue to improve, leading to increased capabilities through competition among them.

My work aims to be model-agnostic regarding foundation model providers. In this project, we studied various proprietary LLMs but also explored using open models like DeepSeek and Llama-3. We found that open models offer significant benefits such as lower costs, guaranteed availability, greater transparency, and flexibility, albeit with slightly lower quality. In the future, we aim to use our proposed discovery process to produce self-improving systems in a closed-loop environment using open models.

# Chapter 5

# Conclusion

## 5.1 Achivement of the Objectives

### 5.1.1 Overview

This project was therefore successful in its key objectives, thus showing that automating significant parts of the scientific research process was possible and feasible. The developed system autonomously produced research ideas, designed and conducted experiments, analyzed the results, and documented the findings in a professional academic format.

### 5.1.2 Idea Generation

In the area of Idea Generation, the system was able to generate a range of new research ideas within a specified scope. This meant that the hypotheses proposed were unique and in line with current scientific trends, thereby fostering innovation within the research domain.

### 5.1.3 Experiment Design and Execution

Regarding Experiment Design and Execution, automated experiments were carried out efficiently; the experiment collected meaningful data in several domains. Being able to run multiple iterations of the experiment contributed to the robustness of findings.

### 5.1.4 Data Analysis and Documentation

For Data Analysis and Documentation, the system correctly analyzed results through standard statistical methods. Detailed reports were produced, complete with visualizations and well-structured write-ups according to academic standards on clarity and presentation.

### 5.1.5 Peer Review Simulation

The project also succeeded in Peer Review Simulation. An automated reviewing mechanism was developed to measure the quality of the research so that only the best ideas would be documented and further improved. This step not only enhanced the credibility of the outputs but also facilitated a systematic approach to quality control in research.

## 5.2 Implications and Recommendations

### 5.2.1 General Implications

The implications of this project are far-reaching, particularly in how research can be conducted and disseminated more efficiently. By automating the research pipeline, the system has the potential to accelerate the pace of scientific discovery and democratize access to research tools. This empowerment can significantly benefit smaller institutions, startups, and independent researchers who may lack the resources for traditional research methodologies.

### 5.2.2 Key Recommendations

However, there are several recommendations for improving the system and ensuring its broader applicability:

- **Scalability:** Future iterations should focus on scaling the system to handle

more complex experiments, larger datasets, and interdisciplinary domains. This could involve integrating more advanced computational resources or optimizing algorithms to run more efficiently. Enhancing scalability will ensure that the system remains relevant across various fields of research.

- **Refinement of Novelty Detection:** The novelty-checking mechanism needs further refinement to improve its ability to detect redundant research ideas, especially in highly specialized fields. Integrating more comprehensive databases and providing real-time access to the latest publications will help address this issue, ensuring that generated ideas are truly innovative.

- **Ethical Review and Safety:** Incorporating a more robust ethical review process is essential to avoid generating potentially harmful or unethical research ideas. This may include additional human oversight or the implementation of more advanced AI-driven ethical review mechanisms. Strengthening ethical safeguards will enhance trust in the system and promote responsible research practices.

- **Collaboration and Integration:** Future versions could benefit from enhanced collaboration features that allow for real-time input from human researchers during the ideation or experimental design phases. This integration would help fine-tune the system's outputs by combining the creativity and expertise of human researchers with the efficiency of automation. Fostering collaboration can lead to richer, more nuanced research outcomes.

## 5.3 Future Scope

### 5.3.1 Development Paths

The potential for extent in the scope of this project is quite huge, with several paths for later development:

41

1. **Cross-Domain Automation:** This framework may be easily adapted for any scientific discipline, such as biology, physics, or social sciences, by incorporating suitable modifications to the phases of data collection and experiment design. It can similarly be used in drug discovery, climate modeling, or in any kind of sociological research that takes too much time for hypothesis generation and subsequent data collection. This aspect of flexibility would make it very useful in a vast variety of studies.

2. **Increased Learning and Adaptation:** Future versions may even include some machine learning approach that makes the system able to "learn" from earlier research cycles. The system, by analysis of past experiments and outcomes, could develop its hypothesis generation and experimental design improvement over time using feedback gained from completed research. Thus, hypotheses as well as experimental approaches are refined further with time.

3. **Interfacing with Physical Laboratories:** Physical experimentation could be the next step of this project with further advancements in robotics and automation. This could be by integrating the system with lab robots, which would automate physical experiments such as material synthesis or genetic analysis to make the pipeline complete from idea generation to lab work. This will bridge the gap between theoretical research and practical application.

4. **Collaborative AI-Human Research Teams:** Developing systems that work collaboratively with human researchers could further enhance the capabilities of the framework. By incorporating user feedback and domain-specific expertise, the system could refine its outputs, creating a synergistic relationship between AI and researchers. This collaboration would leverage the strengths of both human creativity and machine efficiency.

5. **Real-Time Literature Review Integration:** Future developments may include

real-time literature searches and automated synthesis of research papers, so that the system will know the very latest findings in the research field and new ideas will always be based on current scientific knowledge. This will only keep what is suggested really relevant and impactful.

### 5.3.2 Additional Future Directions

Future directions for the framework could include integrating vision capabilities for better plot and figure handling, incorporating human feedback and interaction to refine outputs, and enabling the system to automatically expand the scope of its experiments by pulling in new data and models from the internet, provided this can be done safely. Additionally, the framework could follow up on its best ideas or even perform research directly on its own code in a self-referential manner. Significant portions of the code for this project were written by an AI assistant. Expanding the framework to other scientific domains could further amplify its impact, paving the way for a new era of automated scientific discovery.

For example, by integrating these technologies with cloud robotics and automation in physical lab spaces, the framework could perform experiments in biology, chemistry, and material sciences, provided it can be done safely. Crucially, future work should address reliability and hallucination concerns, potentially through a more in-depth automatic verification of the reported results. This could be achieved by directly linking code and experiments or by determining if an automated verifier can independently reproduce the results.

## 5.4 Personal Reflection

### 5.4.1 Learning Experience

Reflecting on this project, I have gained valuable insights into the power of automation in research and the challenges associated with it. One of the most rewarding

aspects was witnessing how automation could streamline complex workflows, such as experimental design and documentation, which traditionally require significant human effort. This efficiency not only accelerates the research process but also allows researchers to focus more on creative and analytical tasks rather than repetitive procedures.

### 5.4.2 Technical Challenges

The project also provided a deeper understanding of the limitations of current AI technologies, particularly regarding their ability to comprehend complex domain-specific knowledge and manage unexpected errors during automated processes. These challenges were not merely obstacles; they served as valuable learning opportunities that highlighted areas for improvement and optimization. Recognizing these limitations is crucial for developing more robust systems in the future.

### 5.4.3 Ethical Considerations

Additionally, the ethical considerations surrounding the automation of scientific research were particularly eye-opening. Ensuring that generated research remains valuable and safe requires constant vigilance and ethical oversight. This experience underscored the importance of maintaining human judgment in the loop, especially in research fields that can significantly impact society. The need for a balanced approach—where automation enhances human capabilities without replacing critical ethical considerations—became clear throughout this project.

## 5.5 Summary of the key findings

### 5.5.1 Overall Conclusion

This project demonstrated the potential for automating several key stages of the scientific research process. The system was able to autonomously generate novel re-

search ideas, design experiments, execute tests, analyze results, and document findings in an academic format. Additionally, it simulated a peer-review process to ensure the quality of the generated research.

### 5.5.2 Key Points

- **Feasibility of Full Automation:** The research pipeline can be effectively automated, significantly reducing the time and resources required for scientific discovery. This capability enhances productivity and allows researchers to focus on higher-level analytical tasks.

- **Challenges in Computational Efficiency:** While the system functions well for smaller-scale experiments, larger and more complex tasks require further optimization. Addressing these computational efficiency challenges will be crucial for broader application.

- **Potential for Cross-Domain Application:** The framework has broad applicability across various scientific domains. Future work will be needed to customize the system for specific fields, ensuring that it meets the unique requirements of different areas of research.

- **Importance of Ethical Considerations:** Automated systems must be designed with strong ethical safeguards to prevent misuse or the generation of harmful research. This highlights the necessity of maintaining human oversight in automated processes.

### 5.5.3 Final Remarks

The success of this project underscores the significant impact that automation can have on accelerating scientific discovery, particularly for smaller research teams or institutions with limited resources. By providing a foundation for future advance-

ments in automated research, this project aims to make scientific inquiry more efficient, accessible, and innovative.

# References

[1] Gemini Team et al. "Gemini: A Family of Highly Capable Multimodal Models". In: *arXiv e-prints*, arXiv:2312.11805 (Dec. 2023), arXiv:2312.11805. DOI: `10.48550/arXiv.2312.11805`. arXiv: `2312.11805` [`cs.CL`].

[2] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. URL: `https://www.deeplearningbook.org/`.

[3] Aaron Grattafiori et al. "The Llama 3 Herd of Models". In: *arXiv e-prints*, arXiv:2407.21783 (July 2024), arXiv:2407.21783. DOI: `10.48550/arXiv.2407.21783`. arXiv: `2407.21783` [`cs.AI`].

[4] Chris Lu et al. "The AI Scientist: Towards Fully Automated Open-Ended Scientific Discovery". In: *arXiv e-prints*, arXiv:2408.06292 (Aug. 2024), arXiv:2408.06292. DOI: `10.48550/arXiv.2408.06292`. arXiv: `2408.06292` [`cs.AI`].

[5] OpenAI. "GPT-4 Technical Report". In: *arXiv preprint arXiv:2303.08774* (2023). Available online at `https://arxiv.org/abs/2303.08774`.

[6] Adam Paszke, Sam Gross, Francisco Massa, et al. "PyTorch: An Imperative Style, High-Performance Deep Learning Library". In: *Advances in Neural Information Processing Systems (NeurIPS)*. Vol. 32. 2019. URL: `https://papers.nips.cc/paper_files/paper/2019/hash/9015a7baeddc67fa85f0f9ab443eb861-Abstract.html`.

[7] Betty Shea and Mark Schmidt. "Why Line Search When You Can Plane Search? So-friendly Neural Networks Allow Per-iteration Optimization of Learning and Momentum Rates for Every Layer". In: *arXiv* abs/2406.17954 (2024). URL: `https://arxiv.org/abs/2406.17954`.