
IMPROVE: ITERATIVE MODEL PIPELINE REFINEMENT AND OPTIMIZATION LEVERAGING LLM AGENTS

Eric Xue
University of Toronto
e.xue@mail.utoronto.ca

Zeyi Huang
University of Wisconsin - Madison
zeyihuang@cs.wisc.edu

Yuyang Ji
New York University
yj2669@nyu.edu

Haohan Wang
University of Illinois at Urbana-Champaign
haohanw@illinois.edu

ABSTRACT

Computer vision is a critical component in a wide range of real-world applications, including plant monitoring in agriculture and handwriting classification in digital systems. However, developing high-performance computer vision models traditionally demands both machine learning (ML) expertise and domain-specific knowledge, making the process costly, labor-intensive, and inaccessible to many. Large language model (LLM) agents have emerged as a promising solution to automate this workflow, but most existing methods share a common limitation: they attempt to optimize entire pipelines in a single step before evaluation, making it difficult to attribute improvements to specific changes. This lack of granularity leads to unstable optimization and slower convergence, limiting their effectiveness. To address this, we introduce Iterative Refinement, a novel strategy for LLM-driven ML pipeline design inspired by how human ML experts iteratively refine models, focusing on one component at a time rather than making sweeping changes all at once. By systematically updating individual components based on real training feedback, Iterative Refinement improves stability, interpretability, and overall model performance. We implement this strategy in IMPROVE, an end-to-end LLM agent framework for automating and optimizing object classification pipelines. Through extensive evaluations across datasets of varying sizes and domains, including standard benchmarks and Kaggle competition datasets, we demonstrate that Iterative Refinement enables IMPROVE to consistently achieve better performance over existing zero-shot LLM-based approaches. These findings establish Iterative Refinement as an effective new strategy for LLM-driven ML automation and position IMPROVE as an accessible solution for building high-quality computer vision models without requiring ML expertise.

1 Introduction

Computer vision has emerged as a powerful tool for addressing many complex real-world problems, from plant monitoring in agriculture to handwriting classification in digital systems. However, the process of training computer vision models has grown increasingly complex, involving many different steps such as data augmentation, architecture selection, and hyperparameter tuning. As a result, developing high-performing models is labor-intensive and often demands machine learning (ML) expertise, as well as domain-specific knowledge. Each component must be manually calibrated based on prior experience and detailed analysis of training statistics.

To simplify this process, tools such as Weights & Biases have been developed, offering ML practitioners the ability to track, organize, and optimize model training workflows. Similarly, researchers have been investigating approaches, such as hyperparameter optimization [Bergstra and Bengio, 2012, Bergstra et al., 2011, Snoek et al., 2012, Springenberg et al., 2016, Falkner et al., 2018] and neural architecture search [Elsken et al., 2017, Kandasamy et al., 2018], to automate parts of the ML pipeline. While these platforms and methods provide valuable support, they tend to address isolated aspects of the model development process and still require substantial involvement from human experts.

Recent advances in large language models (LLMs) have opened up new possibilities for automation by enabling AI agents to generate code and reason about design choices. Inspired by this capability, researchers have begun leveraging LLM agents to automate the construction and optimization of ML models [Grosnit et al., 2024, Trirat et al., 2024, Hong et al., 2024a, Li et al., 2024, Guo et al., 2024]. Approaches in this emerging field vary widely in terms of planning, reasoning, and evaluation strategies. However, one common feature of most methods is to optimize the entire pipeline simultaneously before each evaluation attempt, rather than refining each component in isolation. While this "all-at-once" approach can help with rapid early-stage improvements, our experiments reveal a key limitation: as model performance approaches its maximum potential, this strategy struggles to make further progress. The inability to attribute performance changes to specific components makes it difficult to refine the pipeline and push improvements beyond a certain threshold.

To address these limitations, we introduce Iterative Refinement, a new strategy for LLM-driven ML system design inspired by how human experts approach model development. Rather than attempting to optimize an entire pipeline at once, experienced ML practitioners typically analyze performance, adjust specific components, and iteratively refine the design based on training feedback. For example, if they determine that data augmentation needs improvement, they will experiment with different augmentation techniques while keeping the training procedure and model architecture unchanged. Only after evaluating its impact will they refine other aspects of the pipeline. Following this intuition, Iterative Refinement systematically updates one component at a time, evaluates its impact, and refines further using real training feedback. This structured approach enables more stable, interpretable, and controlled improvements by isolating the effects of each change and precisely identifying what drives performance gains.

We implement this strategy in IMPROVE (Iterative Model Pipeline Refinement and Optimization leveraging LLM agents), a fully autonomous framework for designing and optimizing image classification pipelines. IMPROVE emulates the workflow of a team of human ML engineers, taking a dataset as input and autonomously constructing, training, and iteratively refining a model through its multi-agent system. It efficiently manages key tasks such as data preprocessing, architecture selection, and hyperparameter tuning.

Our experiments demonstrate that Iterative Refinement enables IMPROVE to generate higher-performing models consistently compared to modifying entire pipelines at once. IMPROVE also greatly outperforms zero-shot LLM-generated training pipelines and achieves near-human-level performance on standard datasets, including CIFAR-10, TinyImageNet, and various Kaggle competition datasets. Notably, it maintains computational efficiency comparable to that of human practitioners, making it both effective and scalable. Our findings validate Iterative Refinement as a practical strategy for LLM-driven ML automation but also establish IMPROVE as an accessible tool for developing state-of-the-art image classification models without requiring deep ML expertise.

Overall, our work introduces a novel strategy for LLM-based ML system design along with a concrete framework that automates the generation and optimization of object classification pipelines. Specifically, our contributions are:

- We introduce Iterative Refinement, a new design approach in LLM agent frameworks inspired by how human experts iteratively refine models, optimizing individual components incrementally for more stable and controlled improvements.
- We design IMPROVE, a novel framework that applies Iterative Refinement to create and optimize object classification pipelines using specialized LLM agents.
- Through extensive experiments on standard and real-world datasets, we demonstrate that Iterative Refinement ensures more consistent performance gains in pipeline optimization, helping IMPROVE to outperform zero-shot LLM-generated pipelines and achieves human expert-level results.

2 Related Works

2.1 Large Language Models and LLM Agents

Large language models (LLMs) are pre-trained on massive text corpora, often with millions to trillions of parameters. Some of the most well-known LLMs today include GPT-3.5 and GPT-4 [OpenAI et al., 2024] from OpenAI, Claude 3.5 from Anthropic, Mixtral from Mistral, and Llama 3 from Meta. While LLMs can be used directly after pre-training, they often undergo additional training stages, such as supervised fine-tuning (SFT) and reinforcement learning from human feedback (RLHF) [Ouyang et al., 2022]. These stages aim to align LLMs with specific behaviors and objectives, resulting in models capable of performing tasks such as general-purpose chatbots, code generation, and information retrieval.

A subset of LLMs specializes in code generation. Codex [Chen et al., 2021], for example, is a GPT model fine-tuned on publicly available GitHub code with a focus on Python programming. Following Codex, Code Llama was introduced,

extending coding capabilities to various programming languages such as Python, C++, Java, PHP, and TypeScript. More recently, larger natural language focused LLMs, including GPT-3.5 and Llama 3-8B, have also demonstrated significant coding proficiency across multiple languages, without needing to be explicitly fine-tuned on code.

LLMs have also been used to create human-like agents capable of executing complex instructions autonomously. Examples of such LLM-based agents include ReAct [Yao et al., 2023], Reflexion [Shinn et al., 2023], and SwiftSage Lin et al. [2023], which have demonstrated effectiveness in complex reasoning and decision-making tasks. Expanding on this, works like AutoGPT [Significant Gravitas, 2024] enabled LLM agents to move beyond reasoning and autonomously perform actions, such as executing code and receiving feedback from outputs. Another line of research explores the collaboration of multiple LLM agents, each with specialized roles [Hong et al., 2024b, Du et al., 2024, Park et al., 2023]. These multi-agent frameworks have shown that LLMs can effectively assume distinct roles [Tseng et al., 2024], and role specialization enhances their ability to retrieve relevant knowledge, outperforming single-agent systems on tasks requiring reasoning and strategic planning [Sreedhar and Chilton, 2024].

2.2 Automated Machine Learning

Automated Machine Learning (AutoML) is a field of research seeking to automate various stages of the ML pipeline, making it more efficient to develop ML models. Most works focus on automating specific parts of the process, such as data preparation, model architecture selection, and hyperparameter optimization [He et al., 2021].

For instance, works like AutoAugment [Cubuk et al., 2019], Faster AutoAugment [Hataya et al., 2019], DADA [Li et al., 2020], and TrivialAugment [Müller and Hutter, 2021] concentrate on automating the data augmentation component of image classification by applying search strategies over predefined search spaces. Another major focus in AutoML is neural architecture search (NAS), which aims to find the most optimal model architecture from a given search space. A prominent example is Elsken et al. [2017], where a hill-climbing algorithm is used to identify the best convolutional neural network (CNN) architecture. More advanced approaches, such as those by Kandasamy et al. [2018], leverage Bayesian optimization and optimal transport to refine this search.

In addition to NAS, hyperparameter optimization (HPO) is a closely related area that seeks to identify the best set of hyperparameters for a fixed architecture. Random search [Bergstra and Bengio, 2012] is a simple yet efficient method that outperforms traditional grid search [Bergstra et al., 2011] by exploring the search space more effectively. To further enhance the search process, many works incorporate Bayesian optimization in HPO [Snoek et al., 2012, Springenberg et al., 2016, Falkner et al., 2018], enabling better use of past information to guide future decisions.

2.3 AutoML with LLMs

With the recent advancement of LLMs, many researchers have started exploring using LLMs to tackle problems in AutoML, as LLMs offer much greater flexibility in the search space over traditional methods. One of the earliest works in this direction is Yu et al. [2023], which utilizes a fine-tuned Generative Pre-Trained Transformer (GPT) model to create new architectures using crossover, mutation, and selection strategies. Later, many works took advantage of the pre-training of LLMs, which grants them an immense amount of implicit knowledge on neural architectures. GENIUS [Zheng et al., 2023] is one such work, where NAS is performed by simply prompting GPT-4 to generate different configurations, which are evaluated iteratively for a pre-determined number of iterations. In a different direction, Hollmann et al. [2024] proposed the CAAFE method that uses LLMs to automate feature engineering by generating semantically meaningful features for tabular datasets. Zhang et al. [2023] worked on a more comprehensive pipeline, attempting to simultaneously tackle many tasks including object detection, object classification and question answering with LLM-based AutoML, by optimizing the pipeline with LLM-predicted training log. Another work, AutoMMLab [Yang et al., 2024], focuses on automating computer vision model training, guiding LLMs to handle data selection, model selection, and hyperparameter tuning sequentially, based on natural language requests.

Recently, researchers have advanced the integration of LLM agents into AutoML, evolving their use from simple code generation tools to comprehensive systems that can automate a wider range of components in the machine learning pipeline. For instance, Agent K [Grosnit et al., 2024] proposes an end-to-end pipeline, where the agent scrapes data from a Kaggle URL, trains a model, and submits results, leveraging LLM reasoning and Bayesian optimization in the process to optimize the pipeline. Similarly, AutoML-Agent [Trirat et al., 2024] employs a multi-agent framework that automates model generation in a multi-stage pipeline, such as data preprocessing, model design, and hyperparameter optimization, and includes multi-stage verification to ensure correctness. Data Interpreter [Hong et al., 2024a] approaches AutoML through dynamic task decomposition, constructing, adjusting, and executing task and action graphs where each node represents a specific sub-task. Other frameworks focus specifically on data science. For example, AutoKaggle [Li et al., 2024] targets tabular data with a similar multi-agent setup, but includes a library of machine learning functions like missing value imputation and one-hot encoding for the agents to use. DS-Agent [Guo et al., 2024] develops models for

text, time series, and tabular data by utilizing case-based reasoning with insights from Kaggle technical reports and code. Although many works in this direction aim to maximize automation, some others focus on human-AI collaboration. For instance, LAMBDA [Sun et al., 2024] is a human-in-the-loop data science agent that employs multi-agent collaboration for code generation, verification, and execution, similar to other methods, but is able to assist users interactively through an intuitive interface.

3 Method

3.1 Overview

Traditional machine learning (ML) pipeline optimization requires deep ML and domain expertise, as practitioners iteratively adjust data processing, model selection, and training configurations based on training performance. While effective, this traditional workflow is labor-intensive, requiring constant supervision from ML experts at every stage of development.

Recent LLM-based automation methods [Grosnit et al., 2024, Trirat et al., 2024, Hong et al., 2024a, Li et al., 2024, Guo et al., 2024] aim to replace this process by autonomously generating ML pipelines. These methods typically incorporate data-driven reasoning and employ optimization techniques to iteratively refine the pipeline for improved performance. However, most existing approaches optimize the entire pipeline as a single entity and treat performance improvements globally all at once, diverging from human workflows, where ML practitioners refine individual components in an iterative manner.

Global optimization strategies present a number of fundamental challenges that hinder its effectiveness. First, modifying multiple pipeline components simultaneously reduces interpretability, making it difficult to attribute performance changes to specific adjustments. Second, due to this lack of interpretability and the vast search space, performance improvements often occur through chance-based trial and error rather than systematic refinement, leading to instability in results. Finally, while these strategies may initially yield rapid improvements in early iterations, they tend to plateau as they struggle to make meaningful refinements at later stages when the overall performance is already high, resulting in an overall slow convergence.

3.2 Iterative Refinement

To address these limitations, we introduce **Iterative Refinement**, a structured optimization strategy inspired by how human ML practitioners systematically adjust individual pipeline components based on empirical feedback. Applying Iterative Refinement in an LLM agent framework requires dividing the target pipeline into distinct components based on the problem domain. Each component must be generated independently yet easily integrated into a complete pipeline. Following the initial training and evaluation of the pipeline, one single component is selected for modification in an attempt to improve while the remaining components remain unchanged. The updated pipeline undergoes training and evaluation, and modifications are retained only if they result in performance improvements. Otherwise, the system reverts to the previous configuration. The process repeats for a fixed number of times, allowing the framework to refine the pipeline progressively.

Compared to existing methods that apply holistic optimization methods, Iterative Refinement offers three key advantages:

- **Improved Interpretability:** By modifying only one component at a time, performance changes can be directly attributed to specific adjustments.
- **Stable Improvement:** The framework does not rely on finding great combinations of parameters and methods based on chance. Instead, it achieves gradual and systematic improvement across the given iterations.
- **Accelerated Convergence:** Targeted refinements prevent redundant modifications and ensure a structured optimization trajectory, leading to faster and more efficient performance gains.

3.3 IMPROVE

To evaluate the Iterative Refinement strategy, we developed **IMPROVE**, a multi-agent end-to-end LLM framework that autonomously generates and optimizes an object classification pipeline. IMPROVE requires only a dataset as input and independently manages the entire process, from pipeline generation to model training, eliminating the need for any human intervention. We provide a diagram of the workflow of IMPROVE in Figure 1.

IMPROVE consists of the following key agents:

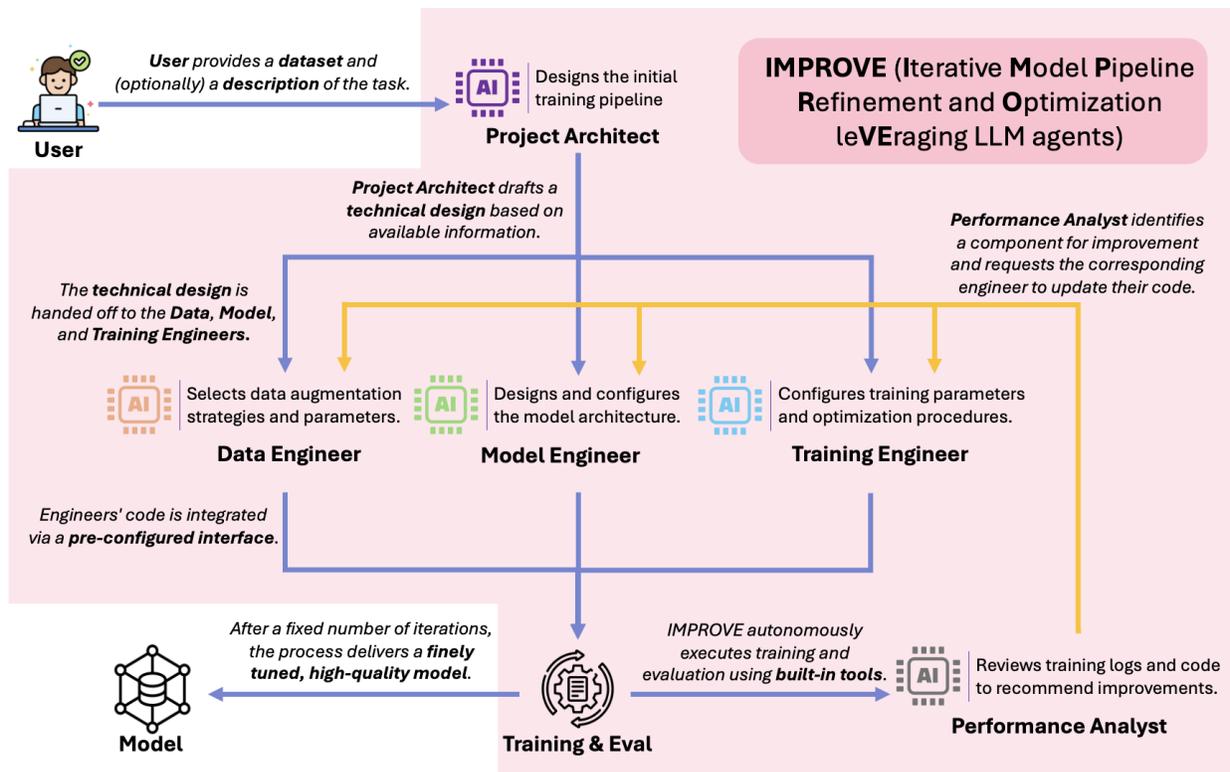


Figure 1: The IMPROVE framework employing Iterative Refinement. Users provide a dataset, and IMPROVE autonomously executes the model development process. The Project Architect designs the pipeline, which is implemented by specialized agents. The Performance Analyst iteratively refines components based on empirical feedback, optimizing the pipeline over multiple iterations. After a set number of iterations, a high-quality model is produced.

- **Project Architect:** This agent is responsible for the initial analysis of the dataset, evaluating its size, structure, and any supplementary information provided. Based on this analysis, the Project Architect generates a comprehensive technical design that serves as the foundation for the entire pipeline.
- **Data Engineer:** This agent handles the implementation of the data processing component, which involves tasks such as data standardization, augmentation, and transformation.
- **Model Engineer:** This agent is tasked with selecting and configuring the model architecture by choosing the most appropriate base model, applying necessary modifications, and incorporating techniques such as regularization if needed.
- **Training Engineer:** This agent is responsible for configuring the model training process, where key tasks include setting the hyperparameters, such as learning rate and batch size, and selecting the best optimization algorithms.
- **Performance Analyst:** This agent is responsible for reviewing pipeline configurations and training logs after each training run is completed. It then identifies one area for improvement and provides targeted feedback to one selected engineer agent.

Each agent in the IMPROVE framework is equipped with knowledge of its role, the team’s goals, and how to collaborate effectively with other agents. Also, each agent operates with domain-specific expertise embedded through prompt engineering. For example, the Model Engineer is aware of effective pre-trained models like ConvNeXt [Liu et al., 2022], while the Data Engineer is familiar with advanced data augmentation techniques such as MixUp [Zhang et al., 2018] and CutMix [Yun et al., 2019]. Although state-of-the-art LLMs implicitly understand these techniques and can implement them when asked to, they rarely apply them to code proactively unless explicitly instructed to do so. Thus, we specify these techniques in the prompts to ensure the agents only select and experiment with highly effective methods when generating the training pipeline.

The framework applies Iterative Refinement to optimize three primary pipeline components that we identified for object classification: **data augmentation**, **model architecture**, and **training procedure**. The development process begins with the **Project Architect**, who reviews the dataset and additional context to produce a technical design. This design is then passed to the **Data Engineer**, **Model Engineer**, and **Training Engineer**, who collaboratively build their respective components according to the plan. Once the data pipeline, model architecture, and training configurations are in place, the system integrates these components and trains the model. The **Performance Analyst** then reviews all past and current pipeline configurations, along with their associated training logs, to identify areas for improvement. It then provides targeted feedback to a selected engineer, who refines their component accordingly. The updated pipeline undergoes retraining and re-evaluation, ensuring that only changes leading to performance gains are retained, while ineffective modifications are logged for future reference but discarded from the active codebase. This iterative cycle continues for a predefined number of iterations, progressively refining the pipeline until the best-performing configuration and corresponding model checkpoint are identified as the final output.

3.4 Implementation

We now highlight some other design choices and innovations in IMPROVE that enhance its effectiveness beyond the Iterative Refinement strategy.

Dataset Information Utilization. A key advantage of the LLM agents framework is its ability to understand any natural language description of the dataset, enabling the system to make more informed decisions. Optionally, the user can supply a brief description of the dataset, detailing aspects such as the image domain (e.g., real-life objects, digits, synthetic images), whether the images are grayscale or colored, their dimensions (same or varying sizes), and their quality (high or low). This information is passed to both the Project Architect and Performance Analyst and is factored into the initial technical design and performance analysis. The impact of including this dataset information is discussed further in Section 4.3.

Autonomous Code Execution. The IMPROVE framework is capable of autonomously generating, executing, and refining model pipeline code over multiple iterations to improve model performance. To achieve this, IMPROVE first needs to ensure that the code produced by each engineer agent is coherent and compatible. This is facilitated by a predefined constant interface that specifies the required input arguments and expected return values for each component. The code-generating agents (i.e., Data Engineer, Model Engineer, and Training Engineer) have flexibility in their implementations but must adhere to these interface guidelines for input and output consistency. To automate the training process, a built-in workflow takes the agents’ responses, and then parses and extracts the code into Python files. The system utilizes the predefined interface to integrate these components into a coherent pipeline and initiates a terminal shell script that executes the full model training process. Training logs and other outputs are captured at the shell level and stored in a text file for analysis.

Unified Pipeline Initialization. At the start of the IMPROVE workflow, the Project Architect agent leverages the strong zero-shot generation capabilities of LLMs to generate a single, coherent pipeline in one step before dividing it into three components. In Section 4.3, we discuss how generating the data, model, and training components separately and then combining them instead can lead to suboptimal initial performance. Using unified pipeline initialization ensures greater efficiency by starting with a well-integrated baseline, reducing the need for extensive corrections and minimizing wasted iterations on poorly constructed initial configurations.

Summarized History. To ensure informed decision-making and prevent redundancy, all previous pipeline configurations and their corresponding training logs are passed to the Performance Analyst. However, passing the full history of code and logs could quickly exceed the LLM’s maximum context length. To mitigate this issue, after each iteration, the system makes a call to the LLM to generate a summary of the key aspects of the data, model, and training code, highlighting information such as the methods and hyperparameters used. Then, instead of passing the entire code from all past iterations, only the summarized configuration along with the training log is provided to the Performance Analyst to guide its reasoning.

4 Experiments

4.1 Experimental Setting

Datasets. We evaluated IMPROVE using two widely recognized classification datasets: CIFAR-10 [Krizhevsky and Hinton, 2009] and TinyImageNet [Deng et al., 2009]. CIFAR-10 contains 60,000 32x32 color images across 10 distinct classes, while TinyImageNet, a subset of ImageNet1K, includes 100,000 images across 200 classes, resized to 64x64 pixels. Both datasets feature commonly seen objects such as vehicles and animals.

To test IMPROVE’s robustness, we also incorporated CIFAR-10-C [Hendrycks and Dietterich, 2019], a variant of CIFAR-10 designed to assess model performance under common corruptions. CIFAR-10-C introduces 15 corruption types, grouped into four categories: noise, blur, weather, and digital distortions. These corruptions are applied at two severity levels (1 and 5), resulting in 30 distinct versions of the dataset. For experimental efficiency, we created two subsets, CIFAR-10-C-1 and CIFAR-10-C-5, where the suffix denotes the severity level. These subsets were generated by uniformly sampling one image from each corrupted dataset for every test image.

Additionally, we evaluated IMPROVE on two smaller datasets from different domains: SVHN and dSprites. SVHN is a digit classification dataset containing real-world images of house numbers, with variations in resolution and background complexity. dSprites is a synthetic dataset of 2D shapes, with latent factors controlling the shapes’ properties. Specifically, we used the dSprites Orientation dataset, where each shape is rotated into one of 40 possible orientations within the $[0, 2\pi]$ range. Both datasets are part of the Visual Task Adaptation Benchmark (VTAB) [Zhai et al., 2020].

To further assess IMPROVE’s versatility across diverse real-world domains, we also evaluated it on four real-world datasets from Kaggle, a leading platform for machine learning competitions. These datasets span a range of sizes and represent various data domains, including:

- **Cassava Leaf Disease Classification** [Mwebaze et al., 2020]: This dataset consists of 21,367 photos of cassava plants in Uganda, mostly taken using inexpensive cameras. The model must identify whether the plant has one of the four diseases: Mosaic Disease, Green Mottle Brown Streak Disease, Bacterial Blight Disease, or no disease, with 5 classes in total.
- **4 Animal Classification** [Lee et al., 2022]: This dataset comprises a total of 3,529 images, categorized into four distinct animal classes: cats, deer, dogs, and horses. Each class represents a diverse range of images capturing various poses, environments, and lighting conditions.
- **Arabic Letters Classification** [Khalil, 2023]: This dataset contains 53,199 images of 65 written Arabic letters, each exhibiting positional variations based on their occurrence within a word with four possible forms: isolated, initial, medial and final. These letters were collected from 82 different users.
- **Kitchenware Classification** [ololo, 2022]: This dataset has 9,367 photos of 6 types of common kitchenware, including cups, glasses, plates, spoons, forks and knives. These photos are taken in households with various lighting and scene conditions.

We used three key criteria to guide our selection of Kaggle datasets to ensure an efficient and fair evaluation. First, we excluded datasets with more than a million images to maintain computational feasibility. Second, we only selected datasets used for public competitions, allowing for a direct comparison between IMPROVE’s performance and that of human ML practitioners. Third, we chose competitions that used top-1 accuracy as the primary evaluation metric, avoiding those that relied on metrics such as the area under the ROC curve, F1 score, or multiclass log loss, to be consistent across our experiments. After applying these filters, these four datasets were among the few that met all the criteria and aligned with the goals of our evaluation.

We always use the provided test set when available. If no test set is provided, we designate the validation set as the test set. In cases where only a single unsplit dataset is available, we manually split the dataset into training and test sets using an 80-20 ratio.

Baseline. For VTAB datasets, namely SVHN and dSprites Orientation, we compare our results with that of Visual Prompt Tuning (VPT) [Jia et al., 2022], a widely recognized method in the domain of fine-tuning. For all four Kaggle datasets, we benchmark IMPROVE against the top existing leaderboard submissions. To further assess the practicality of the IMPROVE framework, we compare the accuracy of its generated model pipelines with that of a straightforward zero-shot prompting approach, which represents the most likely method a non-ML expert would use to generate a model. For this comparison, we provide the LLM with a simple prompt in the format: "Generate code for training a model on a dataset with X classes." The generated script is then executed, and the resulting model is trained and evaluated on the test set manually.

Experimental Setup. We used two commercial LLMs, GPT-4o and o1, developed by OpenAI, for all experiments. The experiments were conducted over three trials, with the average accuracy and their standard deviation reported. For all IMPROVE runs, we performed 20 iterations to balance optimization and experimental efficiency. For the VTAB datasets (SVHN and dSprites Orientation), we instructed the Model Engineer to use the Vision Transformer (ViT), specifically the ViT-B/16 model, to ensure a fair comparison with the VPT paper, which utilized the same model as its baseline.

4.2 Main Results

The results in Table 1 demonstrate that IMPROVE consistently outperforms models generated by zero-shot prompting LLMs using both GPT-4o and o1, confirming its effectiveness as a superior alternative for non-ML experts seeking to train high-performing models. For simpler datasets like CIFAR-10, IMPROVE achieves slightly higher accuracy compared to zero-shot prompting, with the difference being roughly between 3% to 20%. However, its strength becomes more evident on more challenging datasets like TinyImageNet and the two variants of CIFAR-10 with corrupted images. On these datasets, IMPROVE at most achieves a roughly 40% higher accuracy than the zero-shot prompting baseline.

Another notable distinction between the results of IMPROVE and the zero-shot baseline is the significantly lower standard deviation observed with IMPROVE. While zero-shot results occasionally achieve accuracy closer to that of IMPROVE, they are highly inconsistent, even when using the same prompt. For example, the highest recorded standard deviation was approximately $\pm 17\%$ for GPT-4o on CIFAR-10-C-5. This inconsistency means that inexperienced users might achieve decent results at times, but they could just as easily get terrible models. In contrast, this variability is lowered with IMPROVE, as all runs eventually converge to a consistently strong performance.

Table 1: Average classification accuracy for IMPROVE-generated models and zeroshot prompting LLMs on four standard datasets. The best accuracy on each dataset is bolded.

Dataset	IMPROVE (o1)	Zero-shot (o1)	IMPROVE (GPT-4o)	Zero-shot (GPT-4o)
CIFAR-10	0.9825±0.0018	0.7940±0.0287	0.9626±0.0311	0.9290±0.0331
CIFAR-10-C-1	0.9621±0.0026	0.7654±0.0282	0.9476±0.0158	0.5425±0.1007
CIFAR-10-C-5	0.9557±0.0128	0.7662±0.0244	0.9422±0.0038	0.6218±0.1724
TinyImageNet	0.8692±0.0212	0.4630±0.0520	0.7875±0.0169	0.4815±0.1091

Table 2 further illustrates IMPROVE’s effectiveness by comparing it to Visual Prompt Tuning (VPT) [Jia et al., 2022]. IMPROVE not only surpasses VPT but also outperforms the full-parameter fine-tuning baselines reported in the same study. On both datasets, IMPROVE again shows clear improvements over zero-shot prompting, regardless of whether GPT-4o or o1 is used.

Table 2: Average classification accuracy for IMPROVE-generated models and zeroshot prompting LLMs on two VTAB Datasets: SVHN and dSprites Orientation. The table also includes the results from full-parameter fine-tuning (FFT) and Visual Prompt Tuning (VPT) from Jia et al. [2022] for comparison. The best accuracy on each dataset is bolded.

Model	Dataset	IMPROVE	Zeroshot	FFT	VPT
o1	SVHN	0.9779±0.0015	0.9513±0.0027	0.8740	0.7810
	dSprites	0.9667±0.0005	0.5997±0.1407	0.4670	0.4790
GPT-4o	SVHN	0.9695±0.0046	0.9250±0.0307	0.8740	0.7810
	dSprites	0.9523±0.0074	0.6515±0.2199	0.4670	0.4790

In addition to standard datasets, IMPROVE’s performance on Kaggle competition datasets is presented in Table 3. These datasets vary in size, image quality, and domain, offering a more realistic assessment of IMPROVE’s adaptability in real-world scenarios. Across all four datasets, IMPROVE consistently outperforms zero-shot prompting, though its top-1 accuracy remains slightly below the highest-ranking Kaggle leaderboard results. Nevertheless, IMPROVE demonstrates its ability to achieve competitive performance comparable to human ML experts, ranking particularly well on smaller and less complex datasets.

These results were achieved with computational efficiency comparable to that of human experts. Each IMPROVE run was limited to 20 iterations, with more than half of these typically encountering code issues such as undefined variables, wrong package usage, or tensor shape mismatches that cause execution failures. However, these error-containing iterations always terminated quickly, typically within a minute, allowing IMPROVE to proceed to the next iteration with minimal impact on overall efficiency. As a result, fewer than 10 valid iterations are executed and analyzed per run.

This number of valid attempts is comparable to the number of submission attempts made by human practitioners in Kaggle competitions, as shown in Table 3. However, it is important to note that human-submitted Kaggle code is typically tested extensively in local environments on a validation dataset before submission to ensure the code contains no errors and can achieve at least a decent performance. As a result, human practitioners likely conduct far more optimization iterations beyond those reflected in the reported submission counts. Despite these constraints, IMPROVE demonstrates performance close to that of top human ML practitioners, which again highlights its potential as a powerful tool for automated model development in real-world conditions.

Table 3: Average classification accuracy and leaderboard metrics for IMPROVE-generated models and zeroshot prompting LLMs on four Kaggle Datasets. The table also includes the leaderboard rank of the IMPROVE result, the highest accuracy among all Kaggle submissions, and the average submission attempts for the top 5 leaderboard positions.

Model	Dataset	IMPROVE	Zero-shot	Rank	Kaggle Statistics	
					Top Acc.	Top Attempts
o1	Cassava Leaf Disease	0.8931±0.0023	0.8093±0.0163	1591/3900	0.9152	98
	4 Animals	0.9982±0.0015	0.9506±0.0323	1/221	0.9991	12
	Arabic Letters	0.9510±0.0161	0.9162±0.0051	6/177	0.9680	10
	Kitchenware	0.9763±0.0048	0.9044±0.0074	31/115	0.9958	10
GPT-4o	Cassava Leaf Disease	0.8574±0.0275	0.7748±0.0552	2892/3900	0.9152	98
	4 Animals	0.9518±0.0330	0.9196±0.0600	184/221	0.9991	12
	Arabic Letters	0.8403±0.0824	0.5946±0.3264	85/177	0.9680	10
	Kitchenware	0.9793±0.0022	0.8581±0.0956	25/115	0.9958	10

In Table 4, we evaluate the effectiveness of our core technique: Iterative Refinement. To provide a comparison, we developed an alternate version of IMPROVE in which all pipeline components were allowed to be modified simultaneously by their respective agents in every iteration, rather than focusing on improving one component at a time.

The results indicate that this alternative approach led to less coordinated improvements and significantly higher variability in the outcomes. Interestingly, while some individual runs achieved higher performance than IMPROVE, many others performed worse. Observing the trend in Figure 2, we hypothesize that this variability stems from the increased flexibility of modifying multiple components simultaneously. As shown, the approach not using Iterative Refinement expands the search space, offering greater potential to discover an optimal configuration, resulting in fast improvements in the earlier iterations. However, it also introduces a higher risk of failing to converge on a well-performing configuration within the given time constraints. It can be seen that in the later iterations, the version of IMPROVE without Iterative Refinement struggles to find ways to further improve the system, making blind attempts with similar accuracies, whereas IMPROVE with Iterative Refinement can steadily improve its accuracy over time.

Table 4: Average classification accuracy for IMPROVE-generated models, with and without Iterative Refinement (IR), and zeroshot prompting LLMs on CIFAR-10 and TinyImageNet. The best accuracy on each dataset is bolded.

Model	Dataset	IMPROVE	(NO IR)	Zero-shot
o1	CIFAR-10	0.9825±0.0018	0.9579±0.0364	0.7940±0.0287
	TinyImageNet	0.8692±0.0212	0.8339±0.0105	0.4630±0.0520
GPT-4o	CIFAR-10	0.9626±0.0311	0.9610±0.0223	0.9290±0.0331
	TinyImageNet	0.7875±0.0169	0.6202±0.1908	0.4815±0.1091

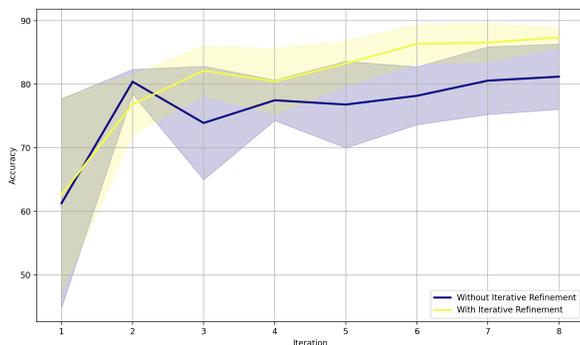


Figure 2: Average accuracy and standard deviation per iteration for a batch of IMPROVE runs on TinyImageNet using o1. The plot compares the performance trajectory of IMPROVE with Iterative Refinement against a variant where Iterative Refinement is removed. Shaded regions around each line indicate the standard deviation. Iterations that resulted in errors, which do not impact accuracy, are omitted.

4.3 Ablation Studies

Unified Pipeline Initialization. As described earlier, IMPROVE employs an initialization strategy known as unified pipeline initialization (UPI). In this approach, the Project Architect generates a complete training pipeline in a single step using zero-shot prompting, which is then broken down into individual components. UPI leverages the LLM’s ability to generate coherent code and recall previously learned configurations, ensuring that the techniques and parameters used across components are well-aligned. As a result, the generated code is more cohesive, leading to a stronger initialization configuration for the model. In contrast, generating each component sequentially - first creating data augmentation code, passing it to the Model Engineer, and then sharing the data and model code with the Training Engineer - often results in less coherent code and reduced model performance. In Table 5, we can see that even though IMPROVE without UPI can still outperform the baseline LLM-generated model, the final accuracy is notably lower than that achieved using UPI. This highlights the importance of a strong initial configuration, which enables IMPROVE to converge to a high-performing model more efficiently.

Table 5: Average classification accuracy for IMPROVE-generated models, with and without unified pipeline initialization (UPI), and zeroshot prompting LLMs on CIFAR-10 and TinyImageNet. The best accuracy on each dataset is bolded.

Model	Dataset	IMPROVE	(NO UPI)	Zero-shot
o1	CIFAR-10	0.9825±0.0018	0.9528±0.0267	0.7940±0.0287
	TinyImageNet	0.8692±0.0212	0.7974±0.0162	0.4630±0.0520
GPT-4o	CIFAR-10	0.9626±0.0311	0.9476±0.0259	0.9290±0.0331
	TinyImageNet	0.7875±0.0169	0.6646±0.0981	0.4815±0.1091

Smaller LLMs. In our experiments, we primarily used two high-performing LLMs, GPT-4o and o1, to build our LLM agents. To test IMPROVE’s robustness with smaller LLMs, we also evaluated its performance using GPT-4o-mini. GPT-4o-mini is presented by OpenAI as the official successor to GPT-3.5, offering improvements in cost, speed, and computational efficiency compared to GPT-3.5 while maintaining a strong performance.

As shown in Table 6, IMPROVE continues to deliver strong results even when using GPT-4o-mini. Despite its smaller size, the model achieves classification accuracies that are significantly higher than those obtained by zero-shot prompting LLMs. This demonstrates IMPROVE’s ability to perform well even with limited budget.

Dataset Information Utilization We also explored how providing additional dataset-specific information influences the design choices and strategies employed by the LLM agents. Specifically, we hypothesize that having access to dataset details would be particularly beneficial for selecting appropriate data augmentation strategies. Our experimental results support this hypothesis, as IMPROVE consistently selected augmentations well-suited to the characteristics of each dataset.

We examined the augmentation strategies IMPROVE employed for each dataset and the reasoning behind them. When training on SVHN, a dataset composed of real-world images of house numbers, IMPROVE selected the ColorJitter augmentation, justifying its choice by stating: "Color jitter (brightness, contrast, saturation, hue) can help in robustifying the model against variations in lighting conditions." This decision is indeed appropriate in this case, as real-world images often contain lighting inconsistencies that can significantly impact model performance.

Conversely, not all augmentations contribute positively to performance. For instance, in the dSprites Orientation dataset, where classification is based on object orientation, applying RandomHorizontalFlip negatively impacts performance by altering class labels and introducing label noise. Although RandomHorizontalFlip was sometimes included in the initial

Table 6: Average classification accuracy for IMPROVE-generated models and zeroshot prompting LLMs on CIFAR-10 and TinyImageNet using GPT-4o mini, with results from o1 and GPT-4o for comparison.

Model	Dataset	IMPROVE	Zero-shot
o1	CIFAR-10	0.9825±0.0018	0.7940±0.0287
	TinyImageNet	0.8692±0.0212	0.4630±0.0520
GPT-4o	CIFAR-10	0.9626±0.0311	0.9290±0.0331
	TinyImageNet	0.7875±0.0169	0.4815±0.1091
GPT-4o-mini	CIFAR-10	0.9590±0.0221	0.8364±0.0493
	TinyImageNet	0.6847±0.1226	0.5781±0.0673

configuration, IMPROVE’s Performance Analyst identified this issue and suggested: "RandomHorizontalFlip may not be useful for orientation classification tasks," and "Orientation-based tasks will not benefit from horizontal flipping; it could even confuse the model."

These findings suggest that IMPROVE is capable of dataset-aware decision-making. Unlike traditional automated approaches that filter through augmentation policies through trial and error, IMPROVE can dynamically adapt its choices based on dataset characteristics and task requirements. By leveraging dataset-specific insights, IMPROVE not only optimizes the augmentation pipeline more effectively but also accelerates model convergence by avoiding ineffective or detrimental transformations.

5 Conclusion

In this paper, we introduced Iterative Refinement, a novel strategy for designing LLM agent-based pipeline optimization systems by decomposing pipelines into independent components and refining each in isolation, enhancing interpretability and consistency while ensuring systematic improvements. To evaluate Iterative Refinement, we implemented IMPROVE, an end-to-end LLM agent framework that autonomously generates image classification models. Beyond Iterative Refinement, IMPROVE integrates several innovative features, including dataset-aware decision-making, automated code execution, and unified pipeline initialization, which together create a robust and adaptive optimization process. Our experiments demonstrated that IMPROVE consistently outperforms models generated through zero-shot LLM prompting, achieves results comparable to top human ML practitioners, and exhibits robustness against corrupted datasets while effectively adapting to diverse image domains. These findings highlight the potential of IMPROVE for real-world applications.

References

- James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.*, 13(null): 281–305, feb 2012. ISSN 1532-4435.
- James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011. URL https://proceedings.neurips.cc/paper_files/paper/2011/file/86e8f7ab32cfd12577bc2619bc635690-Paper.pdf.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021. URL <https://arxiv.org/abs/2107.03374>.
- Ekin D. Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V. Le. Autoaugment: Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255. IEEE, 2009.
- Zhuoyun Du, Chen Qian, Wei Liu, Zihao Xie, Yifei Wang, Yufan Dang, Weize Chen, and Cheng Yang. Multi-agent software development through cross-team collaboration, 2024. URL <https://arxiv.org/abs/2406.08979>.
- Thomas Elsken, Jan-Hendrik Metzen, and Frank Hutter. Simple and efficient architecture search for convolutional neural networks, 2017. URL <https://arxiv.org/abs/1711.04528>.
- Stefan Falkner, Aaron Klein, and Frank Hutter. BOHB: Robust and efficient hyperparameter optimization at scale. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1437–1446. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/falkner18a.html>.
- Antoine Grosnit, Alexandre Maraval, James Doran, Giuseppe Paolo, Albert Thomas, Refinath Shahul Hameed Nabeezath Beevi, Jonas Gonzalez, Khyati Khandelwal, Ignacio Iacobacci, Abdelhakim Benechehab, Hamza Cherkaoui, Youssef Attia El-Hili, Kun Shao, Jianye Hao, Jun Yao, Balazs Kegl, Haitham Bou-Ammar, and

- Jun Wang. Large language models orchestrating structured reasoning achieve kaggle grandmaster level, 2024. URL <https://arxiv.org/abs/2411.03562>.
- Siyuan Guo, Cheng Deng, Ying Wen, Hechang Chen, Yi Chang, and Jun Wang. Ds-agent: Automated data science by empowering large language models with case-based reasoning, 2024. URL <https://arxiv.org/abs/2402.17453>.
- Ryuichiro Hataya, Jan Zdenek, Kazuki Yoshizoe, and Hideki Nakayama. Faster autoaugment: Learning augmentation strategies using backpropagation, 2019. URL <https://arxiv.org/abs/1911.06987>.
- Xin He, Kaiyong Zhao, and Xiaowen Chu. Automl: A survey of the state-of-the-art. *Knowledge-Based Systems*, 212: 106622, January 2021. ISSN 0950-7051. doi: 10.1016/j.knosys.2020.106622. URL <http://dx.doi.org/10.1016/j.knosys.2020.106622>.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.
- Noah Hollmann, Samuel Müller, and Frank Hutter. Large language models for automated data science: introducing caafe for context-aware automated feature engineering. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2024. Curran Associates Inc.
- Sirui Hong, Yizhang Lin, Bang Liu, Bangbang Liu, Binhao Wu, Ceyao Zhang, Chenxing Wei, Danyang Li, Jiaqi Chen, Jiayi Zhang, Jinlin Wang, Li Zhang, Lingyao Zhang, Min Yang, Mingchen Zhuge, Taicheng Guo, Tuo Zhou, Wei Tao, Xiangru Tang, Xiangtao Lu, Xiawu Zheng, Xinbing Liang, Yaying Fei, Yuheng Cheng, Zhibin Gou, Zongze Xu, and Chenglin Wu. Data interpreter: An llm agent for data science, 2024a. URL <https://arxiv.org/abs/2402.18679>.
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. MetaGPT: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2024b. URL <https://openreview.net/forum?id=VtmBAGCN7o>.
- Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *Computer Vision – ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*, page 709–727, Berlin, Heidelberg, 2022. Springer-Verlag. ISBN 978-3-031-19826-7. doi: 10.1007/978-3-031-19827-4_41. URL https://doi.org/10.1007/978-3-031-19827-4_41.
- Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabas Poczos, and Eric P Xing. Neural architecture search with bayesian optimisation and optimal transport. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/f33ba15effa5c10e873bf3842afb46a6-Paper.pdf.
- Youssef Khalil. Arabic letters classification, 2023. URL <https://kaggle.com/competitions/arabic-letters-classification>.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- J H Lee, JHyun Ahn, Sanguk Park, and Seunghyun Jin. 4 animal classification, 2022. URL <https://kaggle.com/competitions/4-animal-classification>.
- Yonggang Li, Guosheng Hu, Yongtao Wang, Timothy Hospedales, Neil M. Robertson, and Yongxin Yang. Dada: Differentiable automatic data augmentation, 2020. URL <https://arxiv.org/abs/2003.03780>.
- Ziming Li, Qianbo Zang, David Ma, Jiawei Guo, Tuney Zheng, Minghao Liu, Xinyao Niu, Yue Wang, Jian Yang, Jiaheng Liu, Wanjun Zhong, Wangchunshu Zhou, Wenhao Huang, and Ge Zhang. Autokaggle: A multi-agent framework for autonomous data science competitions, 2024. URL <https://arxiv.org/abs/2410.20424>.
- Bill Yuchen Lin, Yicheng Fu, Karina Yang, Faeze Brahman, Shiyu Huang, Chandra Bhagavatula, Prithviraj Ammanabrolu, Yejin Choi, and Xiang Ren. Swiftsage: A generative agent with fast and slow thinking for complex interactive tasks, 2023. URL <https://arxiv.org/abs/2305.17390>.
- Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s, 2022. URL <https://arxiv.org/abs/2201.03545>.
- Ernest Mwebaze, Jesse Mostipak, Joyce, Julia Elliott, and Sohier Dane. Cassava leaf disease classification, 2020. URL <https://kaggle.com/competitions/cassava-leaf-disease-classification>.
- Samuel G. Müller and Frank Hutter. Trivialaugment: Tuning-free yet state-of-the-art data augmentation, 2021. URL <https://arxiv.org/abs/2103.10158>.

- ololo. Kitchenware classification, 2022. URL <https://kaggle.com/competitions/kitchenware-classification>.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Aspell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. URL <https://arxiv.org/abs/2203.02155>.
- Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. Generative agents: Interactive simulacra of human behavior, 2023. URL <https://arxiv.org/abs/2304.03442>.
- Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023. URL <https://arxiv.org/abs/2303.11366>.
- Significant Gravitass. AutoGPT, 2024. URL <https://github.com/Significant-Gravitass/AutoGPT>. MIT License.
- Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine learning algorithms. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012. URL https://proceedings.neurips.cc/paper_files/paper/2012/file/05311655a15b75fab86956663e1819cd-Paper.pdf.
- Jost Tobias Springenberg, Aaron Klein, Stefan Falkner, and Frank Hutter. Bayesian optimization with robust bayesian neural networks. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural*

- Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL https://proceedings.neurips.cc/paper_files/paper/2016/file/a96d3afec184766bfeca7a9f989fc7e7-Paper.pdf.
- Karthik Sreedhar and Lydia Chilton. Simulating human strategic behavior: Comparing single and multi-agent llms, 2024. URL <https://arxiv.org/abs/2402.08189>.
- Maojun Sun, Ruijian Han, Binyan Jiang, Houduo Qi, Defeng Sun, Yancheng Yuan, and Jian Huang. Lambda: A large model based data agent, 2024. URL <https://arxiv.org/abs/2407.17535>.
- Patara Trirat, Wonyong Jeong, and Sung Ju Hwang. Automl-agent: A multi-agent llm framework for full-pipeline automl, 2024. URL <https://arxiv.org/abs/2410.02958>.
- Yu-Min Tseng, Yu-Chao Huang, Teng-Yun Hsiao, Wei-Lin Chen, Chao-Wei Huang, Yu Meng, and Yun-Nung Chen. Two tales of persona in llms: A survey of role-playing and personalization, 2024. URL <https://arxiv.org/abs/2406.01171>.
- Zekang Yang, Wang Zeng, Sheng Jin, Chen Qian, Ping Luo, and Wentao Liu. Autommlab: Automatically generating deployable models from language instructions for computer vision tasks, 2024. URL <https://arxiv.org/abs/2402.15351>.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models, 2023. URL <https://arxiv.org/abs/2210.03629>.
- Caiyang Yu, Xianggen Liu, Wentao Feng, Chenwei Tang, and Jiancheng Lv. Gpt-nas: Evolutionary neural architecture search with the generative pre-trained model, 2023. URL <https://arxiv.org/abs/2305.05351>.
- Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features, 2019. URL <https://arxiv.org/abs/1905.04899>.
- Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruysen, Carlos Riquelme, Mario Lucic, Josip Djolonga, Andre Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, Lucas Beyer, Olivier Bachem, Michael Tschannen, Marcin Michalski, Olivier Bousquet, Sylvain Gelly, and Neil Houlsby. A large-scale study of representation learning with the visual task adaptation benchmark, 2020. URL <https://arxiv.org/abs/1910.04867>.
- Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization, 2018. URL <https://arxiv.org/abs/1710.09412>.
- Shujian Zhang, Chengyue Gong, Lemeng Wu, Xingchao Liu, and Mingyuan Zhou. Automl-gpt: Automatic machine learning with gpt, 2023. URL <https://arxiv.org/abs/2305.02499>.
- Mingkai Zheng, Xiu Su, Shan You, Fei Wang, Chen Qian, Chang Xu, and Samuel Albanie. Can gpt-4 perform neural architecture search?, 2023. URL <https://arxiv.org/abs/2304.10970>.