# COMIND: TOWARDS COMMUNITY-DRIVEN AGENTS FOR MACHINE LEARNING ENGINEERING

**Anonymous authors** 

Paper under double-blind review

#### **ABSTRACT**

Large language model (LLM) agents show promise in automating machine learning (ML) engineering. However, existing agents typically operate in isolation on a given research problem, without engaging with the broader research community, where human researchers often gain insights and contribute by sharing knowledge. To bridge this gap, we introduce MLE-Live, a live evaluation framework designed to assess an agent's ability to communicate with and leverage collective knowledge from a simulated Kaggle research community. Building on this framework, we propose CoMind, an multi-agent system designed to actively integrate external knowledge. CoMind employs an iterative parallel exploration mechanism, developing multiple solutions simultaneously to balance exploratory breadth with implementation depth. On 75 past Kaggle competitions within our MLE-Live framework, CoMind achieves a 36% medal rate, establishing a new state of the art. Critically, when deployed in eight live, ongoing competitions, CoMind outperforms 92.6% of human competitors on average, placing in the top 5% on three official leaderboards and the top 1% on one.

# 1 Introduction

The capabilities of large language model (LLM)-based agents are rapidly advancing, showing significant promise in automating complex tasks across domains like software engineering (Jimenez et al., 2023b; Xia et al., 2025), mathematical problem-solving (OpenAI, 2024; Ren et al., 2025; Li et al., 2025), and scientific discovery (Romera-Paredes et al., 2024; Yamada et al., 2025; Sun et al., 2025; Feng et al., 2025). A particularly challenging and impactful frontier for these agents is machine learning engineering (MLE). Automating the multifaceted MLE pipeline, which spans the design, implementation, and rigorous evaluation of high-performance models, remains a critical test of an agent's autonomous reasoning and decision-making abilities.

Recent advances have introduced LLM agents capable of autonomously developing machine learning pipelines for Kaggle-style competitions (Chan et al., 2025). Current approaches have demonstrated a range of techniques, from the ReAct-style reasoning in MLAB (Huang et al., 2024) and the tree-based exploration of AIDE (Jiang et al., 2025), to the skill-specialized multi-agent system of AutoKaggle (Li et al., 2024). Although these systems represent important steps toward automating MLE, they are fundamentally designed to operate in isolation, exploring the solution space individually.

This isolated approach stands in stark contrast to how human experts operate. In real-world data science competitions and research, participants thrive on community knowledge sharing: learning from public discussions, shared code, and collective insights to enhance solution quality and drive innovation (Wuchty et al., 2007). By failing to engage with this dynamic external context, current agents are prone to converging on repetitive strategies and ultimately plateauing in performance. This critical gap motivates our central research question:

How can we evaluate and design research agents that utilize collective knowledge?

To address this question, we introduce **MLE-Live**, a controllable evaluation framework that simulates realistic Kaggle-style research communities with time-stamped public discussions and shared code artifacts that public before competition deadline. This ensure the information access is same as human participant. MLE-Live enables rigorous evaluation of agents' ability to leverage community

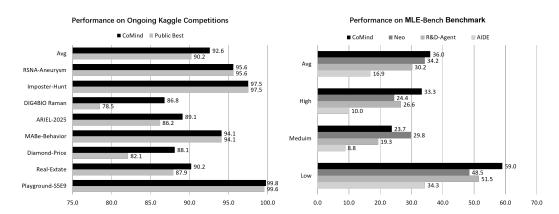


Figure 1: **Left:** CoMind's win rates on eight ongoing Kaggle competitions compared with the public best. **Right:** CoMind achieves state-of-the-art performance on MLE-Bench, measured by the *Any Medal* score.

knowledge in temporally grounded settings, supporting both offline evaluation on past competitions and online evaluation on ongoing competitions.

Building upon this framework, we propose **CoMind**, a multi-agent system designed to systematically incorporate external knowledge and iteratively refine solutions. CoMind's architecture consists of five specialized agent role operating in concert. A central *Coordinator* manages the overall workflow and community interactions. To process external knowledge, an *Analyzer* first summarizes and suggests on improvements and weaknesses for a curated group of solutions, while an *Idea Proposer* brainstorms a diverse pool of ideas and synthesizes novel strategies. These strategies are then passed to multiple parallel *Coding Agents* for implementation and report generation. Finally, a dedicated *Evaluator*, which creates robust scripts for solution assessment and selection. This collaborative process allows CoMind to effectively utile external community knowledge and construct novel solution for the targeted research problem.

We conducted a comprehensive, two-pronged evaluation to assess CoMind's performance in both static and live environments. First, on a static benchmark comprising 75 past Kaggle competitions from MLE-Bench (Chan et al., 2025), CoMind achieved an overall medal rate of 0.36, establishing a new state of the art by significantly outperforming prior leading agents such as Neo and ML-Master (Liu et al., 2025). Second, to validate its real-world practicality, we deployed CoMind in eight ongoing Kaggle competitions (detailed in Figure 1). In this challenging live setting, CoMind proved highly effective, achieving an average rank better than 92.6% of human competitors while placing in the top 5% on three official leaderboards and the top 1% on one. These results demonstrate CoMind's robust effectiveness against contemporary challengers.

In summary, our contributions are:

- MLE-Live: A live evaluation framework simulating community-driven machine learning research with realistic shared discussions and code.
- **CoMind**: A novel agent excelling at collective knowledge utilization and iterative exploration, achieving medal-level performance in real competitions.
- Community-Driven Multiagent Collaboration: An iterative parallel exploration mechanism enabling continuous knowledge accumulation.

#### 2 Related Work

The rise of large language models (LLMs) has sparked a new wave of research into LLM-driven agents, systems that leverage LLMs' reasoning and language capabilities to autonomously perceive, plan, and act within digital or physical environments. Early works such as ReAct (Yao et al., 2023; Schick et al., 2023; Shen et al., 2023; Hong et al., 2023; Boiko et al., 2023) introduced frameworks that transform LLMs into programmable reasoning engines by interleaving natural language reasoning with tool-use actions. Subsequent studies have extended these agents to various domains,

including computer usage (Xie et al., 2024; Zhou et al., 2024) and software development (Wang et al., 2025; Jimenez et al., 2023a).

In parallel, the field of automated machine learning (AutoML) aims to reduce human involvement in building ML pipelines by automating tasks such as model selection, hyperparameter tuning, and architecture search. Early systems like Auto-WEKA (Thornton et al., 2013), HyperBand (Li et al., 2018) and Auto-sklearn (Feurer et al., 2022) used early stopping and Bayesian optimization to search over pipeline configurations, while methods like DARTS (Liu et al., 2019) expanded automation to neural architectures. More recent frameworks such as AutoGluon (Erickson et al., 2020) and FLAML (Wang et al., 2021) emphasize efficiency and ease of use.

Building on these developments, recent efforts have applied LLM-based agents to machine learning engineering (MLE) tasks (Hollmann et al., 2023; Guo et al., 2024; Li et al., 2024; Grosnit et al., 2024; Hong et al., 2024; Chi et al., 2024; Trirat et al., 2024; Huang et al., 2024). However, most evaluations remain constrained to closed-world settings with predefined search spaces, offering limited insight into how these agents perform in open-ended or collaborative ML environments. While some agents (Guo et al., 2024; AI-Researcher, 2025) incorporate basic retrieval tools, these are typically based on simple semantic matching, and robust evaluation methodologies remain underdeveloped.

Meanwhile, several benchmarks have been proposed to evaluate machine learning (ML) engineering capabilities. MLPerf (Mattson et al., 2020) assesses system-level performance, including training speed and energy efficiency. To evaluate end-to-end ML workflows, MLAB (Huang et al., 2024) tests the capabilities of LLM-based agents across 13 ML tasks. MLE-Bench (Chan et al., 2025) and DSBench (Jing et al., 2025) further extends to about 75 Kaggle competitions covering tasks such as preprocessing, modeling, and evaluation. However, these benchmarks typically evaluate agents in isolation, overlooking the collaborative dynamics of real-world ML development. In contrast, our work introduces a framework that simulates community-driven settings, enabling evaluation of agents' ability to engage with and benefit from shared knowledge, while ensuring that resource access remains fair and realistic.

# 3 MLE-LIVE

Existing machine learning benchmarks typically evaluate agents in static, isolated environments. This approach fails to capture the dynamic and collaborative nature of real-world platforms like Kaggle, where progress is driven by community knowledge sharing. Participants constantly learn from shared code, public discussions, and the iterative work of others, making these community interactions a decisive factor in developing top-tier solutions.

To bridge this gap, we introduce **MLE-Live**, a live evaluation framework that extends the widely-used MLE-Bench (Chan et al., 2025). The core innovation of MLE-Live is its simulation of community interactions, providing agents with a time-stamped stream of discussions and code artifacts that mirrors the natural flow of public knowledge during a competition.

Each competition environment in MLE-Live includes the following components: (i) Task description: The background, specifications, evaluation metrics, and data structure, scraped directly from the original Kaggle competition. (ii) Competition dataset: A cleaned train-test split of the official data. When necessary, this includes reconstructed test sets to account for data that is no longer public. (iii) Submission grader: An evaluation script that precisely mimics Kaggle's official scoring mechanism. (iv) Leaderboard: A snapshot of the final public leaderboard. (v) Community artifacts: A curated set of discussions and code notebooks that were **published before the competition deadline**. These artifacts are enriched with valuable metadata (e.g., vote counts, public scores, author tiers) to signal quality and are accompanied by any public datasets or models they reference, creating a self-contained and realistic research environment.

MLE-Live aggregates a substantial dataset of 12,951 discussions and 15,733 kernels from 75 Kaggle competitions. To ensure fairness and eliminate post-hoc data leakage, it strictly includes only resources available prior to competition deadlines, forcing agents to operate under the same information constraints as human participants. This approach offers numerous benefits for robust evaluation: it grounds agents in diverse, objectively-graded ML problems from Kaggle, while the controlled information scope allows for a fair assessment of their retrieval and reasoning abilities. These features enhance reproducibility and enable consistent, longitudinal comparisons between different agents.

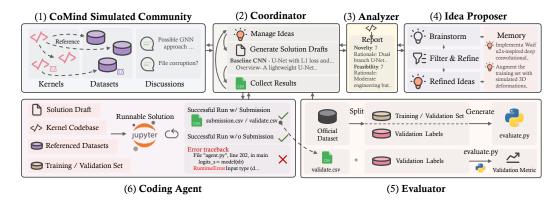


Figure 2: Overview of CoMind. Specialized agents (Coordinator, Analyzer, Idea Proposer, Coding Agent, Evaluator) interact with a simulated community of kernels, datasets, and discussions.

# 4 CoMIND

We propose **CoMind**, a community-augmented large language model (LLM) agent designed to automate machine learning (ML) engineering in an iterative, collaborative setting. Figure 2 is an overview of CoMind workflows.

#### 4.1 COMMUNITY SIMULATION

CoMind's effectiveness stems from simulating the collaborative dynamics that drive breakthrough performance in competitive ML environments. Unlike isolated automated ML systems, CoMind replicates how Kaggle participants leverage community knowledge: drawing insights from discussions, adapting public notebooks and datasets, and contributing discoveries back to the collective knowledge pool.

The simulated community is represented as  $(\mathcal{K}_t, \mathcal{D}_t, \mathcal{T}_t)$  at iteration t, where  $\mathcal{K}_t$  contains all kernels with evaluation metrics,  $\mathcal{D}_t$  includes published datasets and model checkpoints, and  $\mathcal{T}_t$  captures the dependency relationships between resources. CoMind initializes a high-quality community  $(\mathcal{K}_0, \mathcal{D}_0, \mathcal{T}_0)$  by fetching  $k_{\text{kernel}}$  top-performing kernels and  $k_{\text{discussion}}$  most popular discussions from Kaggle, along with all referenced datasets and models. The system constructs a dependency graph  $\mathcal{T}_0 = (V, E)$  where vertices represent kernels or datasets and edges capture resource dependencies.

This dependency structure enables CoMind to systematically trace solution construction, identify influential artifacts, and prioritize resources that drive performance improvements. The graph facilitates intelligent ensemble strategies by combining complementary approaches while avoiding redundant components.

CoMind operates as an active community participant, iteratively analyzing promising kernels, generating novel solutions, conducting experiments, and contributing successful results back to the community. Each iteration produces new artifacts: enhanced kernels, augmented datasets, or ensemble checkpoints, that expand the community knowledge base with associated performance metrics.

Through this continuous cycle of exploration and contribution, CoMind simulates the collaborative dynamics of competitive ML development, where collective intelligence progressively advances performance frontiers at automated scale and speed.

### 4.2 MULTI-AGENT SYSTEM

CoMind orchestrates machine learning experimentation through a coordinated multi-agent system. Specialized agents collaborate in distinct roles, mirroring the division of expertise in human research teams across ideation, implementation, and evaluation. The workflow is an iterative loop managed by the Coordinator, which delegates tasks to the other agents.

Coordinator The Coordinator serves as CoMind's central orchestration hub. Its primary responsibilities are managing the workflow, interfacing with the community environment, and allocating resources. At the start of each iteration t, the Coordinator initiates the process by strategically sampling promising code notebooks (kernels)  $\mathcal{K}_t'$  and relevant datasets  $\mathcal{D}_t'$  from the community. This focused sampling directs the system's attention toward high-potential areas. After receiving refined ideas from the Idea Proposer, the Coordinator translates them into concrete solution drafts  $\mathcal{S}_t$ , which are comprehensive blueprints detailing model architecture, feature engineering, and training procedures. It then instantiates multiple Coding Agents in parallel, assigning each a distinct draft and all referenced resources. Upon completion, the Coordinator aggregates the results and publishes successful solutions back to the community, advancing the environment state for the next iteration.

Analyzer The Analyzer is responsible for distilling raw community artifacts into structured, actionable intelligence. It receives the sampled kernels and discussions from the *Coordinator* and performs a deep analysis across four key dimensions: novelty, feasibility, effectiveness, and efficiency. For each artifact, it generates a 0-10 score on these metrics, accompanied by qualitative explanations of successful patterns, emerging trends, or potential pitfalls. The output is a set of structured analytical reports  $\mathcal{R}_t$ , which serve as the primary input for the *Idea Proposer*.

Idea Proposer The Idea Proposer functions as CoMind's creative engine, tasked with generating novel solution concepts. It uses the analytical reports  $\mathcal{R}_t$  from the Analyzer and its own persistent memory of historical ideas  $\mathcal{I}_t^*$  to ensure that new concepts are both innovative and informed by past results. The ideation process follows three phases: (1) Brainstorming: Generating a wide array of diverse ideas, prioritizing creativity and exploration. (2) Filtering: Ranking these ideas based on feasibility, potential for improvement, and alignment with the analytical reports. Only the most promising subset of ideas  $\mathcal{I}_t$  is selected. (3) Memory Integration: Updating its knowledge base with the newly generated ideas ( $\mathcal{I}_{t+1}^* = \mathcal{I}_t^* \cup \mathcal{I}_t$ ), allowing for increasingly sophisticated strategies over time. The final output, a filtered set of high-potential ideas  $\mathcal{I}_t$ , is sent back to the Coordinator to be developed into full solution drafts.

**Coding Agent** The *Coding Agent* is the implementation workhorse, responsible for converting the abstract solution drafts from the *Coordinator* into executable code. Following an iterative, ReActstyle approach, it conducts trial-and-error experiments using the training and validation data provided by the *Evaluator*. To maximize efficiency, the agent maintains a persistent Jupyter Notebook session to eliminate data reloading overhead and employs a monitor LLM to track execution and terminate failed runs immediately. This iterative process of coding, debugging, and optimization continues until a viable solution is produced or a time budget is exhausted.

Evaluator The Evaluator ensures objective, standardized, and reproducible assessment across all experiments, mirroring official Kaggle protocols. It first partitions the public dataset D into a training set  $D^*$  and a validation set with inputs  $V_x$  and ground-truth labels  $V_y$ . Crucially, only  $D^*$  and  $V_x$  are accessible to the Coding Agents, preserving the integrity of the validation process. When a Coding Agent submits predictions  $V_{\hat{y}}$ , the Evaluator computes the performance score using the official competition metric  $\varphi(V_{\hat{y}}, V_y)$ . It maintains a global leaderboard of all experimental runs, enabling CoMind to reliably track progress and make informed decisions about which solutions to prioritize and publish.

#### 5 BENCHMARK EVALUATION

#### 5.1 SETUP

**Task Selection.** Based on MLE-Live evaluation framework, we evaluate our agent on 75 Kaggle competitions on MLE-Bench. Using the MLE-Live framework, CoMind has access to shared discussions and public kernels published on the competition websites before the competition deadline. Since the MLE-bench test set may be constructed from Kaggle's official public training set, and publicly available datasets or model checkpoints may have been trained on this portion of the data, we restricted CoMind's access to public datasets to minimize potential data contamination. It can only view code published by other contestants.

Table 1: Any Medal (%) scores on 75 MLE-Bench competitions. CoMind achieves state-of-the-art results across difficulty levels. Best results in each column are bolded. Baseline numbers are taken from the official MLE-Bench leaderboard.

| Agent                        | Low (%) | Medium (%) | High (%) | All (%) |
|------------------------------|---------|------------|----------|---------|
| CoMind o4-mini               | 59.09   | 23.68      | 33.33    | 36.00   |
| Neo multi-agent              | 48.48   | 29.82      | 24.44    | 34.22   |
| R&D-Agent o3 + GPT-4.1       | 51.52   | 19.30      | 26.67    | 30.22   |
| ML-Master deepseek-r1        | 48.50   | 20.20      | 24.40    | 29.30   |
| R&D-Agent o1-preview         | 48.18   | 8.95       | 18.67    | 22.40   |
| AIDE o1-preview              | 34.30   | 8.80       | 10.00    | 16.90   |
| AIDE gpt-4o                  | 19.00   | 3.20       | 5.60     | 8.60    |
| AIDE claude-3-5-sonnet       | 19.40   | 2.60       | 2.30     | 7.50    |
| OpenHands gpt-4o             | 11.50   | 2.20       | 1.90     | 5.10    |
| AIDE llama-3.1-405b-instruct | 8.30    | 1.20       | 0.00     | 3.10    |
| MLAB gpt-4o                  | 4.20    | 0.00       | 0.00     | 1.30    |

To validate CoMind under realistic conditions, we further evaluate CoMind on eight ongoing Kaggle competitions. These competitions span diverse domains, including tabular learning, text regression, image classification and video recognition. Rather than approximating the official scoring locally, we directly submit CoMind's generated submission.csv files to the Kaggle platform, so that all reported ranks reflect genuine, live leaderboard positions.

**Implementation Details.** CoMind employs o4-mini-2025-04-16 (OpenAI, 2025) as its backend LLM. We limit the hardware constraint of each run to 32 vCPUs and a single A6000 GPU. Each competition is evaluated in separate containers with a maximum of 24 hours to produce the final submission file. Every single code execution session is limited to 5 hour. Each Coder is limited to a maximum of 30 steps. The number of parallel agents is set to 4.

During code generation, agents are provided with the test set inputs (without labels) and prompted to generate a submission.csv file. The submission is then evaluated by a grader that compares the predicted labels with the ground truth. Following the setting of MLE-Bench, to avoid potential overfitting, test set labels and the competition leaderboard are strictly withheld from the agent's accessible environment. Instead, each agent must rely solely on a self-constructed "runtime test set", a held-out split from the original training data, for code evaluation and performance estimation.

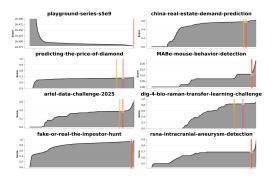
**Metrics.** Following the evaluation metrics in MLE-Bench, we measure the performance of Co-Mind by **Any Medal**, the percentage of competitions where the agent earns a gold, silver, or bronze medal.

Baselines. We compare CoMind against the MLE-Bench leaderboard<sup>1</sup> including open-sourced systems like **R&D-Agent** (Yang et al., 2025), a dual-agent framework (Researcher/Developer) that explores multiple solution branches and merges promising ideas into improved pipelines; **ML-Master** (Liu et al., 2025), which integrates exploration and reasoning via a selectively scoped memory that aggregates insights from parallel trajectories; **AIDE** (Jiang et al., 2025), a purposebuilt tree-search scaffold that iteratively drafts, debugs, and benchmarks code for Kaggle-style tasks; **OpenHands** (Wang et al., 2025), a general-purpose CodeAct-based scaffold that executes code and calls tools in a sandboxed environment; **MLAB** (Huang et al., 2024), referring to the ResearchAgent scaffold from MLAgentBench, a general tool-calling/plan—act baseline; and **Neo** (https://heyneo.so/), a close-sourced multi-agent system for autonomous ML engineering.

#### 5.2 RESULTS

Table 1 compares CoMind with baseline methods on 75 MLE-Bench competitions. CoMind achieves state-of-the-art performance with an *Any Medal* rate of 36.00%, significantly outper-

<sup>1</sup>https://github.com/openai/mle-bench



| Competition       | Rank | Teams | Тор % |
|-------------------|------|-------|-------|
| Playground S5E9   | 4    | 1966  | 0.2%  |
| China Real Estate | 43   | 437   | 9.8%  |
| Diamond Price     | 8    | 67    | 11.9% |
| MABe Behavior     | 3    | 51    | 5.9%  |
| ARIEL 2025        | 90   | 827   | 10.9% |
| DIG4BIO Raman TL  | 22   | 167   | 13.2% |
| Impostor Hunt     | 26   | 1037  | 2.5%  |
| RSNA Aneurysm     | 35   | 788   | 4.4%  |

Figure 3: **Left:** Score distributions across participants in eight ongoing Kaggle competitions. Each curve shows the relationship between leaderboard rank (x-axis, inverted) and competition score (y-axis). Vertical lines indicate CoMind's position (red) and public best performance (yellow). **Right:** Results on eight ongoing Kaggle competitions. Reported are leaderboard rank, total teams, and percentile rank (Top %, where lower means better standing).

forming open-source competitors such as R&D-Agent (submitted on 2025-08-15) and surpassing the closed-source multi-agent system Neo. Appendix C provides a detailed case study on denoising-dirty-documents.

On the eight evaluated ongoing competitions, CoMind ranked top 7.35% on average and improved the best public kernel on 5 competitions. Details including authentic scores and win rates per task are provided in Figure 3. These authentic results demonstrate CoMind's capability to tackle a variety of problem domains and achieve competitive performance in live, evolving ML workflows.

### 6 ABLATION STUDY

#### 6.1 SETUP

**Task Selection.** To evaluate the impact of introducing public resources, we conducted an ablation study on 20 competitions from MLE-Bench-Lite based on MLE-Live. These tasks span across various categories, including image classification/generation, text classification/generation, image regression, audio classification, and tabular analysis.

**Baselines.** We compared CoMind against the following baselines. For consistency, all baselines use the same backend model as CoMind:

- AIDE+Code. To enable the use of publicly available code (e.g., Kaggle kernels), we extend AIDE with access to one public kernel per draft node, which is selected by highest community votes. AIDE+Code augments the prompt with both the task description and the selected kernel alongside the tree summarization.
- AIDE+RAG. We further equip AIDE with a retrieval-augmented generation (RAG) mechanism. Before generating code, the agent retrieves the titles of the top 10 voted discussions and kernels. The LLM selects the most relevant ones, receives a summarization, and then proposes its plan and implementation. For debugging or refinement, it can optionally re-query documents. Retrieval is based on cosine similarity between query and candidate document embeddings, using Multilingual E5 Text Embeddings (Wang et al., 2024).
- CoMind w/o  $\mathcal{R}$ .  $\mathcal{R}$  denotes all public resources. In this variant, CoMind operates without access to any external community resources. It starts with an empty community and relies solely on its own generation history to propose candidate ideas and assemble solution drafts.

**Metrics.** Following the evaluation metrics in prior research (Chan et al., 2025), the relative capability of generating high-quality solution compared with human is measured by:

- Above Median: Indicates whether the submission outperforms at least 50% of competitors on the leaderboard.
- Win Rate: The percentage of competitors whose final scores are lower than the agent's score. If the agent fails to produce a valid submission, the Win Rate is 0.

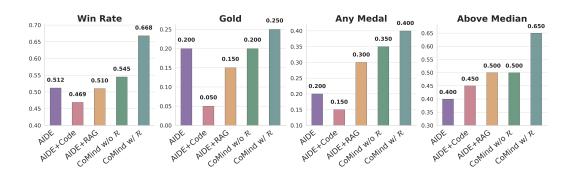


Figure 4: **Performance of CoMind and other baselines on 20 competitions from MLE-Bench- Lite.** *Valid Submission* is the ratio of submissions meeting format requirements and validation criteria. *Win Rate* is the percentage of human competitors outperformed by the agent. *Any Medal*, is the proportion of competitions where the agent earned Gold, Silver or Bronze medals. *Above Median* is the fraction of competitions where the agent's score strictly exceeded the median human competitor.

Table 2: Average win rate of CoMind and other baselines across task categories on 20 competitions from MLE-Bench-Lite. # of Tasks refers to the number of competitions in the corresponding category. CoMind consistently outperforms baselines across most domains.

| Category             | # of Tasks | CoMind | AIDE+Code | AIDE+RAG | AIDE  |
|----------------------|------------|--------|-----------|----------|-------|
| Image Classification | 8          | 0.597  | 0.459     | 0.434    | 0.525 |
| Text Classification  | 3          | 0.740  | 0.157     | 0.338    | 0.61  |
| Audio Classification | 1          | 0.901  | 0.272     | 0.259    | 0.271 |
| Seq2Seq              | 2          | 0.408  | 0.503     | 0.550    | 0.228 |
| Tabular              | 4          | 0.664  | 0.673     | 0.688    | 0.483 |
| Image To Image       | 1          | 0.988  | 0.932     | 0.617    | 0.568 |
| Image Regression     | 1          | 0.992  | 0.342     | 0.992    | 0.992 |
| All                  | 20         | 0.668  | 0.469     | 0.510    | 0.512 |

- Medals: Medals are assigned based on the agent's score relative to Kaggle leaderboard thresholds for gold, silver, and bronze medals.
- Any Medal: The percentage of competitions in which the agent earns any medal.

**Implementation Setup.** All agents use o4-mini-2025-04-16 as their backend. Based on the settings of our main experiment, the hardware constraint is further limited to 4 vCPUs and 5 hours per competition. Each execution session is limited to 1 hour. Access to public datasets are restricted. In accordance with baselines, CoMind has access to 10 top-voted discusions and kernels.

#### 6.2 RESULTS

Figure 4 shows the results. Our key findings are as follows: (i) CoMind consistently outperforms all baselines across every metric. (ii) Among the AIDE variants, AIDE+RAG outperforms AIDE+Code, and both surpass the original AIDE on most metrics, demonstrating the benefits of integrating community knowledge. CoMind further exceeds these approaches, highlighting the effectiveness of its deeper and more strategic community-aware exploration. (iii) Removing CoMind's resource access causes a significant drop in valid submission rates and other metrics, showing that strategic access to public resources helps CoMind balance extending established methods for reliability with exploring novel approaches.

# 7 ANALYTICAL EXPERIMENTS

For analytical experiments, we adopt the same setup as the ablation study and evaluate model performance across multiple dimensions, including task categories, win rate over time, and code complexity.



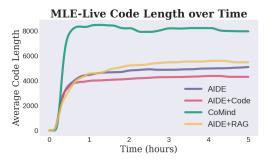


Figure 5: Win rate over time. CoMind sustains improvement while baselines plateau.

Figure 6: Code complexity over time. CoMind generates longer, richer solutions than baselines.

**Task Categories** Table 2 reports the average ranks across seven task categories. CoMind outperforms all baselines in Image Classification, Text Classification, Audio Classification, and Image-to-Image tasks, highlighting its strong adaptability. We manually inspect the tasks where CoMind underperformed and find that the issues are often related to the use of large models or datasets. For example, in Seq2Seq tasks, CoMind explores complex fine-tuning strategies for large language models which often fail to complete within the one-hour runtime constraint.

Win Rate Over Time Figure 5 shows the evolution of average win rate over time. AIDE quickly produces concise, functional solutions, leading to a rapid rise in performance during the first hour. In contrast, CoMind spends more time on debugging and exploration early on, resulting in a slower initial improvement. However, after the first two hours, AIDE's performance plateaus, while CoMind continues to improve through iterative refinement and deeper exploration, ultimately surpassing AIDE and achieving higher-quality solutions.

**Code Complexity** Regarding code complexity, Figure 6 illustrates the average code length during the entire competition. CoMind consistently generates significantly longer and more complex code, while other baselines begin with simpler implementations and introduce only incremental modifications. Appendix A offers a comparative analysis across code complexity metrics and task categories. Notably, CoMind's solutions for Image Regression and Audio Classification are nearly twice as long as those of other baselines. Additionally, solutions from CoMind are, on average, 55.4% longer than those produced by AIDE.

#### 8 CONCLUSION

We introduced MLE-Live, the first framework to evaluate ML agents in community-driven settings, simulating the collaborative dynamics that are essential to real-world progress in Kaggle competitions and beyond. Building upon this benchmark, we proposed CoMind, a community-augmented LLM agent that iteratively selects and synthesizes ideas, implements solutions, and shares reports within a simulated ecosystem. Our results demonstrate that CoMind not only achieves state-of-theart performance on retrospective MLE-Bench tasks but also attains medal-level standings in live Kaggle competitions.

**Limitations and Future Work.** While our current experiments focus on Kaggle-style ML tasks, the MLE-Live framework can be extended to broader domains, such as scientific discovery, openended coding, or robotics, enabling research agents to contribute meaningfully across diverse fields.

# REFERENCES

- AI-Researcher. Ai-researcher: Fully-automated scientific discovery with llm agents, 2025. URL https://github.com/HKUDS/AI-Researcher. Accessed: 2025-05-15.
- Daniil A. Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research with large language models. *Nature*, 624:570 578, 2023. URL https://api.semanticscholar.org/CorpusID:266432059.
- Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, Aleksander Madry, and Lilian Weng. MLEbench: Evaluating machine learning agents on machine learning engineering. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=6s5uXNWGIh.
- Yizhou Chi, Yizhang Lin, Sirui Hong, Duyi Pan, Yaying Fei, Guanghao Mei, Bangbang Liu, Tianqi Pang, Jacky Kwok, Ceyao Zhang, Bangbang Liu, and Chenglin Wu. Sela: Tree-search enhanced llm agents for automated machine learning. *ArXiv*, abs/2410.17238, 2024. URL https://api.semanticscholar.org/CorpusID:273507330.
- Nick Erickson, Jonas Mueller, Alexander Shirkov, Hang Zhang, Pedro Larroy, Mu Li, and Alexander Smola. Autogluon-tabular: Robust and accurate automl for structured data. *arXiv preprint arXiv*:2003.06505, 2020.
- Shengyu Feng, Weiwei Sun, Shanda Li, Ameet Talwalkar, and Yiming Yang. A comprehensive evaluation of contemporary ml-based solvers for combinatorial optimization. *ArXiv*, abs/2505.16952, 2025. URL https://arxiv.org/abs/2505.16952.
- Matthias Feurer, Katharina Eggensperger, Stefan Falkner, Marius Lindauer, and Frank Hutter. Autosklearn 2.0: Hands-free automl via meta-learning. *Journal of Machine Learning Research*, 23 (261):1–61, 2022. URL http://jmlr.org/papers/v23/21-0992.html.
- Antoine Grosnit, Alexandre Max 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, Balázs Kégl, Haitham Bou-Ammar, and Jun Wang. Large language models orchestrating structured reasoning achieve kaggle grandmaster level. *ArXiv*, abs/2411.03562, 2024. URL https://api.semanticscholar.org/CorpusID:273850235.
- Siyuan Guo, Cheng Deng, Ying Wen, Hechang Chen, Yi Chang, and Jun Wang. Dsagent: Automated data science by empowering large language models with case-based reasoning. *ArXiv*, abs/2402.17453, 2024. URL https://api.semanticscholar.org/CorpusID:268033675.
- Noah Hollmann, Samuel G. Müller, and Frank Hutter. Large language models for automated data science: Introducing caafe for context-aware automated feature engineering. In *Neural Information Processing Systems*, 2023. URL https://api.semanticscholar.org/CorpusID:258547322.
- Sirui Hong, Xiawu Zheng, Jonathan P. Chen, Yuheng Cheng, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zi Hen Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, and Chenglin Wu. Metagpt: Meta programming for multi-agent collaborative framework. *ArXiv*, abs/2308.00352, 2023. URL https://api.semanticscholar.org/CorpusID:260351380.
- Sirui Hong, Yizhang Lin, Bangbang Liu, Binhao Wu, Danyang Li, Jiaqi Chen, Jiayi Zhang, Jinlin Wang, Lingyao Zhang, Mingchen Zhuge, Taicheng Guo, Tuo Zhou, Wei Tao, Wenyi Wang, Xiangru Tang, Xiangtao Lu, Xinbing Liang, Yaying Fei, Yuheng Cheng, Zhibin Gou, Zongze Xu, Chenglin Wu, Li Zhang, Min Yang, and Xiawu Zheng. Data interpreter: An Ilm agent for data science. *ArXiv*, abs/2402.18679, 2024. URL https://api.semanticscholar.org/CorpusID:268063292.
- Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. MLAgentbench: Evaluating language agents on machine learning experimentation. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=1Fs1LvjYQW.

Zhengyao Jiang, Dominik Schmidt, Dhruv Srikanth, Dixing Xu, Ian Kaplan, Deniss Jacenko, and Yuxiang Wu. Aide: Ai-driven exploration in the space of code, 2025. URL https://arxiv.org/abs/2502.13138.

- Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *ArXiv*, abs/2310.06770, 2023a. URL https://api.semanticscholar.org/CorpusID: 263829697.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023b.
- Liqiang Jing, Zhehui Huang, Xiaoyang Wang, Wenlin Yao, Wenhao Yu, Kaixin Ma, Hongming Zhang, Xinya Du, and Dong Yu. DSBench: How far are data science agents from becoming data science experts? In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=DSsSPrORZJ.
- Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. *Journal of Machine Learning Research*, 18(185):1–52, 2018. URL http://jmlr.org/papers/v18/16-558.html.
- Shanda Li, Tanya Marwah, Junhong Shen, Weiwei Sun, Andrej Risteski, Yiming Yang, and Ameet Talwalkar. Codepde: An inference framework for llm-driven pde solver generation. *arXiv* preprint arXiv:2505.08783, 2025.
- Ziming Li, Qianbo Zang, David Ma, Jiawei Guo, Tuney Zheng, Minghao Liu, Xinyao Niu, Yue Wang, Jian Yang, Jiaheng Liu, et al. Autokaggle: A multi-agent framework for autonomous data science competitions. *arXiv preprint arXiv:2410.20424*, 2024.
- Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: Differentiable architecture search. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=S1eYHoC5FX.
- Zexi Liu, Yuzhu Cai, Xinyu Zhu, Yujie Zheng, Runkun Chen, Ying Wen, Yanfeng Wang, E Weinan, and Siheng Chen. Ml-master: Towards ai-for-ai via integration of exploration and reasoning. *ArXiv*, abs/2506.16499, 2025. URL https://api.semanticscholar.org/CorpusID:279465426.
- Peter Mattson, Christine Cheng, Gregory Diamos, Cody Coleman, Paulius Micikevicius, David Patterson, Hanlin Tang, Gu-Yeon Wei, Peter Bailis, Victor Bittorf, et al. Mlperf training benchmark. *Proceedings of Machine Learning and Systems*, 2:336–349, 2020.
- OpenAI. Learning to reason with llms, 2024. URL https://openai.com/index/learning-to-reason-with-llms/.
- OpenAI. Introducing openai o3 and o4-mini, 2025. URL https://openai.com/index/introducing-o3-and-o4-mini/.
- Z. Z. Ren, Zhihong Shao, Junxiao Song, Huajian Xin, Haocheng Wang, Wanjia Zhao, Liyue Zhang, Zhe Fu, Qihao Zhu, Dejian Yang, Z. F. Wu, Zhibin Gou, Shirong Ma, Hongxuan Tang, Yuxuan Liu, Wenjun Gao, Daya Guo, and Chong Ruan. Deepseek-prover-v2: Advancing formal mathematical reasoning via reinforcement learning for subgoal decomposition, 2025. URL https://arxiv.org/abs/2504.21801.
- Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools, 2023. URL https://arxiv.org/abs/2302.04761.

Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face, 2023. URL https://arxiv.org/abs/2303.17580.

- Weiwei Sun, Shengyu Feng, Shanda Li, and Yiming Yang. Co-bench: Benchmarking language model agents in algorithm search for combinatorial optimization. *ArXiv*, abs/2504.04310, 2025. URL https://arxiv.org/abs/2504.04310.
- Chris Thornton, Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Auto-weka: Combined selection and hyperparameter optimization of classification algorithms. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 847–855, 2013.
- Patara Trirat, Wonyong Jeong, and Sung Ju Hwang. Automl-agent: A multi-agent llm framework for full-pipeline automl. *ArXiv*, abs/2410.02958, 2024. URL https://api.semanticscholar.org/CorpusID:273162376.
- Chi Wang, Qingyun Wu, Markus Weimer, and Erkang Zhu. Flaml: A fast and lightweight automl library. *Proceedings of Machine Learning and Systems*, 3:434–447, 2021.
- Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Multilingual e5 text embeddings: A technical report. *arXiv* preprint arXiv:2402.05672, 2024.
- Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert Brennan, Hao Peng, Heng Ji, and Graham Neubig. Openhands: An open platform for AI software developers as generalist agents. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=OJd3ayDDoF.
- Stefan Wuchty, Benjamin F. Jones, and Brian Uzzi. The increasing dominance of teams in production of knowledge. *Science*, 316:1036 1039, 2007. URL https://api.semanticscholar.org/CorpusID:260992737.
- Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. Demystifying llm-based software engineering agents. *Proceedings of the ACM on Software Engineering*, 2(FSE):801–824, 2025.
- Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. OSWorld: Benchmarking multimodal agents for open-ended tasks in real computer environments. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL https://openreview.net/forum?id=tN61DTr4Ed.
- Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist-v2: Workshop-level automated scientific discovery via agentic tree search. *arXiv preprint arXiv:2504.08066*, 2025.
- Xu Yang, Xiao Yang, Shikai Fang, Bowen Xian, Yuante Li, Jian Wang, Minrui Xu, Haoran Pan, Xinpeng Hong, Weiqing Liu, et al. R&d-agent: Automating data-driven ai solution building through llm-powered automated research, development, and evolution. *arXiv preprint arXiv:2505.14738*, 2025.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=WE vluYUL-X.
- Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=oKn9c6ytLx.

# A ADDITIONAL ANALYSIS ON CODE COMPLEXITY

In this section, we provide a comprehensive analysis of the generated code using a broad set of software complexity and quality metrics, beyond mere line counts. Specifically, we report the following indicators: Cyclomatic Complexity (CC), Pylint score, Halstead Metrics: Volume, Difficulty, Effort, Source Lines of Code (SLOC), Number of Comment Lines and Code Length.

Table 3: Code complexity and quality metrics (Cyclomatic Complexity, Pylint score, Halstead metrics, SLOC, etc.) across task categories. CoMind produces more complex solutions compared to baselines.

| Category             | Metric        | CoMind   | AIDE    | AIDE+RAG | AIDE+Code |
|----------------------|---------------|----------|---------|----------|-----------|
| Image Classification | CC            | 1.68     | 1.59    | 1.93     | 1.29      |
|                      | Pylint Score  | 7.43     | 9.06    | 8.90     | 8.92      |
|                      | Volume        | 330.88   | 143.26  | 84.20    | 175.88    |
|                      | Difficulty    | 4.95     | 2.90    | 2.32     | 2.59      |
|                      | Effort        | 1960.22  | 507.06  | 286.31   | 725.59    |
|                      | SLOC          | 198.25   | 133.50  | 120.88   | 115.71    |
|                      | Comment Lines | 15.62    | 12.88   | 13.75    | 14.43     |
|                      | Code Length   | 7638.40  | 4624.30 | 4701.30  | 5192.10   |
|                      | CC            | 3.58     | 4.28    | 2.00     | 0.00      |
|                      | Pylint Score  | 8.82     | 9.09    | 8.89     | 9.26      |
|                      | Volume        | 286.38   | 384.07  | 47.68    | 29.25     |
| T4 Cl: 64:           | Difficulty    | 3.76     | 3.94    | 1.25     | 1.31      |
| Text Classification  | Effort        | 1183.11  | 2332.22 | 61.56    | 35.16     |
|                      | SLOC          | 181.67   | 133.00  | 141.00   | 69.50     |
|                      | Comment Lines | 14.67    | 15.33   | 14.00    | 13.50     |
|                      | Code Length   | 6974.70  | 3094.50 | 5920.50  | 5629.30   |
|                      | CC            | 2.00     | 0.00    | 0.00     | 0.00      |
|                      | Pylint Score  | 7.92     | 9.11    | 9.49     | 8.86      |
|                      | Volume        | 718.63   | 244.20  | 115.95   | 227.48    |
| A 4: - C1 : £ 4:     | Difficulty    | 7.39     | 6.46    | 3.19     | 6.38      |
| Audio Classification | Effort        | 5308.07  | 1577.11 | 369.58   | 1451.30   |
|                      | SLOC          | 256.00   | 82.00   | 92.00    | 72.00     |
|                      | Comment Lines | 20.00    | 11.00   | 16.00    | 16.00     |
|                      | Code Length   | 9449.00  | 3508.00 | 4151.00  | 3352.00   |
|                      | CC            | 4.38     | 2.25    | 22.33    | 15.75     |
|                      | Pylint Score  | 8.58     | 9.04    | 9.14     | 8.51      |
|                      | Volume        | 492.55   | 52.33   | 390.46   | 324.00    |
| Seq2Seq              | Difficulty    | 3.87     | 2.14    | 5.26     | 3.68      |
| seqzseq              | Effort        | 1935.02  | 140.58  | 2083.84  | 1686.74   |
|                      | SLOC          | 184.50   | 63.50   | 222.50   | 147.50    |
|                      | Comment Lines | 22.50    | 13.00   | 23.00    | 19.50     |
|                      | Code Length   | 6925.50  | 5649.50 | 8357.50  | 2728.50   |
|                      | CC            | 2.78     | 1.62    | 2.38     | 0.25      |
|                      | Pylint Score  | 8.65     | 8.96    | 8.87     | 9.31      |
| Tabular              | Volume        | 1264.61  | 856.12  | 815.29   | 435.46    |
|                      | Difficulty    | 7.37     | 4.83    | 6.05     | 3.69      |
|                      | Effort        | 10808.93 | 6163.62 | 5564.22  | 2001.06   |
|                      | SLOC          | 218.75   | 139.75  | 147.50   | 93.50     |
|                      | Comment Lines | 18.25    | 14.75   | 15.25    | 10.50     |
|                      | Code Length   | 8570.00  | 3534.00 | 6064.00  | 5759.80   |
|                      | CC            | 1.72     | 2.00    | 3.00     | 1.88      |
|                      | Pylint Score  | 8.43     | 6.25    | 6.64     | 7.74      |
|                      | Volume        | 1298.11  | 1481.62 | 414.59   | 431.08    |
| Image to Image       | Difficulty    | 9.68     | 6.73    | 3.94     | 3.79      |
|                      | Effort        | 12565.66 | 9967.24 | 1633.22  | 1631.93   |
|                      | SLOC          | 228.00   | 175.00  | 121.00   | 128.00    |
|                      | Comment Lines |          | 8.00    | 23.00    |           |

| Category         | Metric        | CoMind   | AIDE    | AIDE+RAG | AIDE+Code |
|------------------|---------------|----------|---------|----------|-----------|
|                  | Code Length   | 8800.00  | 5231.00 | 4815.00  | 6671.00   |
| Image Regression | CC            | 1.68     | 2.00    | 2.40     | 2.00      |
|                  | Pylint Score  | 8.62     | 8.75    | 8.80     | 8.89      |
|                  | Volume        | 1310.92  | 241.08  | 70.32    | 72.00     |
|                  | Difficulty    | 8.75     | 3.88    | 2.18     | 2.73      |
|                  | Effort        | 11466.58 | 934.17  | 153.43   | 196.36    |
|                  | SLOC          | 267.00   | 145.00  | 116.00   | 133.00    |
|                  | Comment Lines | 36.00    | 15.00   | 12.00    | 12.00     |
|                  | Code Length   | 10991.00 | 4841.00 | 4655.00  | 5614.00   |

# B PROMPTS AND RESPONSES FOR COMIND

This section provides some examples of prompts and responses in CoMind, including **Coordinator**, **Analyzer**, **Idea Proposer**, **Coding Agent** and **Evaluator**.

#### B.1 COORDINATOR

#### **Prompt for Solution Draft Synthesis**

**Introduction** You are an expert machine learning researcher preparing for the Kaggle competition described below.

**Task Description**  $\{description \ of \ the \ specified \ task\}$ 

**Ideas** {*entries in the idea pool*}

**Reports** {entries in the report pool}

**Public Pipelines** {all public pipelines extracted before}

Goals

- 1. Carefully read the reports provided above.
- 2. Based on the ideas and reports, propose {num\_pipes} promising self-contained pipelines that are likely to perform well.
- 3. The Public pipelines section contains top-ranked public pipelines during the competition. Use them as reference to polish your pipelines.
- 4. Each pipeline should not overlap with others. Your proposed pipelines should include **one baseline pipeline that uses well-known methods but is robust and relatively easy to implement**. You should reinforce public pipelines and previous pipelines based on their reports (if provided).
- Ensure that each pipeline can be trained within 2 hours on a single A6000 with 48GB memory.
- 6. Read the **submission format** requirements in the task description carefully. The format requirement is possible to be different from the training dataset. **THIS IS EXTREMELY IMPORTANT**. Mention in the pipeline descriptions and be sure to include the code that handles the input and output.
- 7. DO NOT USE tensorflow, use pytorch instead

# Response Template for Solution Draft Synthesis

Submit Pipelines Descriptions and codes of pipelines, separated each pipeline by ===SEP-ARATOR=== mark. For each pipeline, attach code that captures its essential. You must include the code in public pipelines that handles input and output, and if there are parts of the public pipelines that are similar to the current pipeline, you should include them as well.

# B.2 ANALYZER

# **Prompt for Strategy Distilation**

**Introduction** You are an expert machine learning researcher preparing for the Kaggle competition described below.

**Task Description** { description of the specified task }

**Goals** These are top-ranked public scripts during the competition. Your job is to:

- 1. Carefully read the following scripts.
- 2. For each script, if it's self-contained, i.e., including model architecture (if there's a model), training strategies, evaluation, etc., then summarize its pipeline.
- 3. If the pipeline contains technical details, such as extensive feature engineering, hyperparameter tuning, etc., then list them in full detail.

4. Select a representative code segment for each pipeline. You must include dataset reading / submission generation parts. If task-specific details such as feature engineering are included, the code segment should contain them as well.

**Public Kernels** { contents of public kernels}

# Response Template of Strategy Distillation of Public Kernels

**Pipelines** Description of each strategy, separated by ===SEPARATOR=== mark. For each strategy, follow this format:

- Pipeline: A full detailed description of the pipeline. All input/output format, hyperparameters, training settings, model architectures, feature engineering, validation metric, and any other relevant information should be included. Do not omit any feature engineering details.
- Code abstract: A representative code segments that captures the essence (including in-put/output) and novelty of the pipeline. You MUST go through all the publicly available code and include the parts that generate the submission file. Contain task-specific engineering details. Mark the remainder as ellipses.

# **Prompt for Strategy Distillation of Public Discussions**

**Introduction** You are an expert machine learning researcher preparing for the Kaggle competition described below.

**Task Description** { description of the specified task}

**Goals** These are top-voted public discussions during the competition. Your job is to:

**Public Discussions** {contents of public discussions}

- 1. Carefully read the following discussions.
- 2. For each discussion, you should decompose it into critical, novel and inspiring ideas that have potential to win this competition.

### Response Template of Strategy Distillation of Public Discussions

**Ideas** required format: python list of strings, each element is a description of an idea extracted from the discussions. e.g. ['idea 1', 'idea 2'].

#### B.3 IDEA PROPOSER

#### **Prompt for Brainstorm**

**Introduction** You are an expert machine learning researcher preparing for the Kaggle competition described below.

**Task Description** { description of the specified task}

**Goals** I already have a list of ideas that partially explore how to approach this competition. Your job is to:

- Think creatively and construct at least 4 alternative and highly novel solution
  paths that are likely to perform well, especially if combined with careful experimentation.
- 2. Each solution path can be a strategy, pipeline, or method that combines multiple techniques. Try to make them as different as possible from the existing "ideas" list.
- 3. After describing each full solution path, **break it down into individual minimal ideas**-these should be the smallest units of implementation (e.g., "use LightGBM for baseline", "normalize input features", "apply stratified K-fold CV")
- 4. Ensure these ideas do not substantially duplicate items already in "ideas".

5. Refer to the "Reports" section for the latest updates and suggestions on the ideas and previous pipelines.

**Ideas** {entries in the idea pool}

**Reports** {entries in the report pool}

**Public Pipelines** {all public pipelines extracted before}

**Instructions** Format your output like this (one line, one idea):

# **Response Template**

```
{your understanding of the task and explanation of your approaches}
===SOLUTION_PATH_1===
{description of this approach}
- {minimal idea 1}
- {minimal idea 2}
- {minimal idea 3}
- ...
===SOLUTION_PATH_2===
...
===SOLUTION_PATH_3===
...
```

Be ambitious but realistic - many ideas can later be tested on a small subset of the data. Focus on novelty, diversity, and decomposability. Ready? Start.

### **Prompt for Idea Filtering and Reconstruction**

**Introduction** You are a machine learning expert. After carefully searching the relevant literature, you have come up with a list of ideas to implement. However, this idea list has some issues:

- Some ideas are too similar and should be merged into one.
- Some ideas are overlapping, you should rephrase and decouple them.
- You should discard ideas that are irrelevant to the final performance, such as error visualization, etc.

You should refer to the Reports section and Public Pipelines section for previous implemented pipelines. Please decompose, merge, and reconstruct the ideas listed below.

**Ideas** {*entries of the idea pool*}

**Reports** {*entries of the report pool*}

**Public Pipelines** {all public pipelines extracted before}

# Response Template of Idea Filtering and Reconstruction

**Ideas** required format: Python list of strings, each element is a description of an idea. e.g. ['idea 1', 'idea 2'].

# Prompt for Coding Agent Report Compilation

Please summarize the results and submit a comprehensive report.

# **Response Template for Coding Agent Report Compilation**

**pipeline** A detailed description of the pipeline that generated the best results. All hyperparameters, training settings, model architectures, feature engineering, validation metric, and any other relevant information should be included. Describe potential improvements and future work.

**summary** A comprehensive evaluation of each individual component of the pipeline. For each component, summarize in the following format:

=== {name of the component} ===

**Novelty**: 0-10 (0: trivial, 10: clearly novel - major differences from existing well-known methods)

{your rationale}

**Feasibility**: 0-10 (0: almost impossible to implement and require extensive engineering, 10: Easy to implement)

{your rationale}

**Effectiveness**: 0-10 (0: minimal performance improvement, 10: very strong performance, significantly outperform most baselines)

{your rationale}

**Efficiency**: 0-10 (0: very slow, over-dependent on CPU and hard to produce meaningful results within the time limit, 10: high utilization of GPU)

{your rationale}

**Confidence**: 0-10 (0: no emprical results, not sure whether the evaluation is correct, 10: fully verified on large scale with abundant results)

#### B.4 EVALUATOR

#### Prompts for Dataset Splitting and evaluate.py

You are an experienced machine learning engineer. Please generate two self-contained Python code for local evaluation of a Kaggle agent. Your code should be robust, reusable, accept command-line arguments and print necessary information.

#### Background

- Kaggle competitions usually provide labels only for the training set. To evaluate an agent locally, we need to split the training set into a training and validation split.
- The validation set must hide its labels from the agent. The agent only sees the training set (with labels) and the validation inputs (without labels).
- The hidden validation labels will be stored separately and used only for offline evaluation.
- Importantly: ./public must never contain validation labels. Validation labels are saved only in ./private.

**Kaggle Competition Description** { description of the specified task}

**Data Preview** {schema of the input file structure}

**Deliverables** Please generate two scripts (both in Python 3, runnable from the command line):

#### 1) split\_dataset.py

Goal: Split the original training data into 90% training and 10% validation. Store validation inputs (without labels) in ./public, and validation labels in ./private. The training set (with labels) and original test set must remain in ./public, preserving the original structure as closely as possible. The structure of validation inputs should also match the test set. Generate a sample validate submission validate\_sample\_submission.csv under ./public. All original data (training and test) are visible in {path to the input directory}.

**Example**: If the original data is structured as:

```
- kaggle_evaluation/ (official evaluation tool provided by Kaggle)
- __init__.py
- ...
- train.csv
```

```
- train/
- test.csv
- test/
- sample_submission.csv
You should split the dataset into:
(./public/)
 kaggle\_evaluation/ (official evaluation tool provided by Kaggle) (unchanged, soft
    links)
    ___init_
- train.csv (this contains 90% of the training data)
- train/ (this contains 90% of the training data, keep unchanged data as soft links)
- test.csv (unchanged, soft link)
- test/ (unchanged, soft link)
- sample_submission.csv (unchanged, soft link)
- validate.csv (this contains 10% of the training data with labels withheld)
 validate/ (soft links)
- validate_sample_submission.csv (a sample submission file for validation set)
(./private/)
 validate.csv (labels of validation set)
```

If the training data contains zip files, you should extract them to the public directory before splitting the dataset. You should always print the directory structure after the split. Do not extract files to the original directory and keep it unchanged.

If the training data contains multiple classes, you should use **stratified sampling**. You should strictly follow the evaluation metric mentioned in the task description and ensure the validation set is representative of the overall class distribution. Never write validation labels into ./public.

Your code will be executed by command line as follows:

DO NOT store the training and test files in other folders such as ./public\_iTIMESTAMP;, the ./public folder will be exposed to later code agent. Make sure the ./public directory has similar structure with the original data folder.

#### 2) evaluate.py

Goal: Evaluate the agent's predictions on validation set against the hidden ground truth (./private/...). Output evaluation results (json format) to console and write ./private/e-val\_report.json.

It will be executed by command line as follows:

```
'`'bash
python evaluate.py --public\_dir ./public --private\_dir ./private --pred <path to the
   validation submission file>
'''
```

We will pass the path to the sample validation submission file as the argument to your evaluate.py script. It typically produces low scores.

The script should generate in the following json format at ./private/eval\_report.json:

```
"score": A float number represents the evaluation score on the validation set. Do
    not omit this field. If the evaluation is unsuccessful or the predictions are
    invalid, this field should be set to null,
    "success": A boolean value indicates whether the evaluation was successful or not,
    "message": A string provides additional information about the evaluation result.
    Leave it an empty string if the predictions are valid and evaluation is
    successful. Otherwise provide necessary details on why it failed.
```

Do not raise any error or exception. If the evaluation is unsuccessful, you should set the score to null and provide a detailed explanation in the message field.

Now, let's write these two scripts step by step. Your should first generate split\_dataset.py. We will execute the code by command line as mentioned above. You should correct the code in case of any issues. You should always generate full, self-contained code. No part of the code should be omitted.

# Respond in the following format: '''current\_file This should be either split\_dataset.py or evaluate.py. Leave this as None if both are generated and functioned. This indicates the current file you are editing. '''explanation You explanation on the workflow of your code. '''python The full content of the current file. Leave this as None if both are generated and functioned. '''

#### B.5 CODING AGENT

# **Prompts for Coding Agent Iterative Implementation**

**Introduction** You're an expert Kaggle competitor tasked with implementing a pipeline into Python code. You can modify the details (training parameters, feature engineering, model selection, etc.), but do not change overall architecture of this pipeline. The goal is to **obtain best score** on this competition.

**Task Description** { description of the specified task}

**Pipeline** { description of the solution draft to implement }

**Data Overview** { schema of the input file structure} Follow the pipeline description and the code abstract to implement it. All the input files are visible in ../input folder, this folder typically contains the competition data and external resouces, including public datasets, models and outputs of other kernels. DO NOT USE /kaggle/input paths in your code. USE ../input instead.

file structure:

```
- input/ (../input)
  - competition_id/ # the official competition dataset
  - alice/dataset1/ # other public datasets
  - alice/kernel1/ # referenced kernels
- working/
  - agent.ipynb # the notebook you will be working on (./agent.ipynb)
  - other files
```

You will develop the pipeline based on this codebase. Any output files of the codebase, such as csvs, checkpoints, etc., are visible in ./, which is also your current working directory. {Description of Selected Codebase}

You should note that checkpoints generated by this codebase is store in ./ other than ../input. You must load the checkpoint file under the ./ directory for ensemble prediction.

Your code must produce a submission at ./submission.csv, this is EXTREMELY IMPORTANT. Before generating the submission, you must print the value of the evaluation metric computed on a hold-out validation set. You can use custom evaluation functions during training, but the final metric MUST FOLLOW THE EVALUATION SECTION IN THE TASK DESCRIPTION on a validation set. If other kernels with submission.csv are provided in the input folder, you can ensemble them before generating your own submission. This is important because we will pick your best code based on this metric. You are allowed to load the checkpoints of other models. Do not contain any absolute paths in your code. Time limit per run is 2 hours. Your code will be killed if timeout.

Your code will be executed on a single A6000 GPU. Use large batchsizes to maximize the gpu utilization. If the code segment is provided in this prompt, you should follow the input/output structure. You are allowed to install any packages you need or inspect the workspace (e.g., print file contents, check folder structure). Always use gpu for acceleration. DO NOT USE ABSOLUTE PATHS IN YOUR CODE.

The workspace will be maintained across iterations. That is, if your first iteration code produces a checkpoint, you can load it in the second iteration. You can ensemble submissions generated by yourself and other kernels. You should generate model checkpoints for future

loading. If you load the external submissions successfully but failed to merge them with your own predictions, you should print the headers of the external submission and your own predictions and check if the ids are aligned. All the external submissions are valid. Your predictions should be in the same format as them.

To evaluate your submission locally. You should also generate a submission file on the validation set. All the validation data are typically structured similarly to the test data. An external grader will be used to evaluate your validation submission. That is to say, you should generate TWO submission files: one is for the validation set and the other is for the test set. Generate two submission files in the same code cell.

You are allowed to install any packages by running 'pip install ¡package\_name¿' in your script. Your installation will take effect in the NEXT cell.

A persistent Jupyter Notebook session is maintained. Your proposed code cell will be directly appended to the notebook and executed. You should separate data loading, training and evaluation in different cells. Now, please propose THE FIRST CELL of your code (not your full code!) using the following format:

```
The explanation of your first cell. You should describe the desired execution time and
    output of this cell. Explain how to interpret the execution output.
</goal>
The content of this cell. Do not wrap the code in a markdown block. Your code will be
    appended to the notebook, which is stored at ./agent.ipynb. Your code must print
    necessary information after each milestone.
<validation_submission>
The name of the submission file for the validation set. e.g. validate\_submission.csv.
    If your current code cell does not produce two submission files, leave this as
    None.
</validation submission>
<submission>
The name of the submission file for the test set. e.g. submission.csv. This submission
    should be ready for Kaggle submission. If your current code cell does not produce
    two submission files, leave this as None.
</submission>
```

The validation\_submission tag and the submission tag should must be both empty or both non-empty.

## **Prompt for Execution Monitor**

You are an AI assistant monitoring code execution. Your task is to analyze the current execution output and decide whether the code should continue running.

Code being executed:

{code to analyze}

Goal: {execution target of this code}

Runtime Information:

- Current runtime: {code execution time elapsed}
- Maximum runtime: {maximum execution time}
- Remaining time: {remaining execution time}

Current Output:

{current output of this code cell}

1125 1126

1080

1081

1082

1083

1084

1085

1087

1088

1089

1090

1091

1092

1093

1094 1095

1096

1099

1100

1101 1102

1103

1104

11051106

1107

1108

1109

1110

1111

1112 1113

1114 1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1128

1129

1130

1131

1132

1133

Consider these factors:

1126 Consider these

- 1. Is the loss exploding (becoming very large or NaN)?
- 2. Is the loss decreasing normally over time?
- 3. Are there any error messages indicating failure?
- 4. Does the output suggest normal training/execution progress?
- 5. Based on current progress and remaining time, is it possible to complete within the time limit?

# Respond in the following format:

<explanation>
Your rationale for the action. Describe the current progres

Your rationale for the action. Describe the current progress, your estimated remaining time, and explain why you think the execution should continue or stop. DO NOT GIVE SUGGESTIONS ON BUG FIXES.

# **Prompt for Consequent Code Revisions**

The execution takes {execution\_time} seconds and ends with the following output: {truncated output}

Execution completed successfully. You should keep updating your code (e.g., try different hyperparameters, augmentations, model architectures) after you have made successful submission. Your best submission will be recorded.

```
Now, respond in the following format:
<validation_submission>
The name of the submission file for the validation set. e.g. validate_submission.csv.
    If your current code cell does not produce a submission file on the validation set
      leave this as None.
</validation_submission>
<submission>
The name of the submission file for the test set. e.g. submission.csv. This submission
    should be ready for Kaggle submission. If your current code cell does not produce
    a submission file on the test set, leave this as None.
</submission>
<goal>Describe the goal and how to inspect the output of your next code cell</goal>
The content of your next code cell. Following the previous format, do not wrap your
    code within markdown code marks. You should keep updating your code (e.g., try
    different hyperparameters, augmentations, model architectures) even after you have
     made successful submission. Always evaluate your submission and print the metric
    on a validation set.
```

The validation\_submission tag and the submission tag should must be both empty or both non-empty.

# C CASE STUDY: DENOISING DIRTY DOCUMENTS

#### C.1 DATASET PREPARATION

1188

1189 1190

1191 1192

1193

1194 1195

1196

Besides the task description and datasets prepared in MLE-Bench, MLE-Live collects 59 public kernels and 19 discussions which are available on Kaggle and are posted before the competition ends.

#### C.1.1 EXAMPLE OF PUBLIC KERNEL

```
1197
1198 <sup>1</sup>
        A simple feed-forward neural network that denoises one pixel at a time
1199
1200 4
        import numpy as np
        import theano
1201 5
        import theano.tensor as T
1202 <sup>6</sup>
1203 <sup>7</sup>
        import cv2
        import os
1204 _{9}^{\circ}
       import itertools
1205<sub>10</sub>
1206 11
       theano.config.floatX = 'float32'
1207<sup>12</sup>
1208 <sup>13</sup>
        def load_image(path):
1209 14
1209 15
           return cv2.imread(path, cv2.IMREAD_GRAYSCALE)
1210 <sub>16</sub>
        def feature_matrix(img):
            """Converts a grayscale image to a feature matrix
1211 17
1212<sup>18</sup>
1213 19
           The output value has shape (<number of pixels>, <number of features>)
1214 20 21
           # select all the pixels in a square around the target pixel as
1215
               features
1216 22
           window = (5, 5)
           nbrs = [cv2.getRectSubPix(img, window, (y, x)).ravel()
1217 <sup>23</sup>
                   for x, y in itertools.product(range(img.shape[0]), range(img.
1218 <sup>24</sup>
                       shape[1]))]
1219 <sub>25</sub>
1220 <sub>26</sub>
           # add some more possibly relevant numbers as features
1221 27
           median5 = cv2.medianBlur(img, 5).ravel()
           median25 = cv2.medianBlur(img, 25).ravel()
1222 28
1223 <sup>29</sup>
           grad = np.abs(cv2.Sobel(img, cv2.CV_16S, 1, 1, ksize=3).ravel())
1224 30 31
           div = np.abs(cv2.Sobel(img, cv2.CV_16S, 2, 2, ksize=3).ravel())
1225 <sub>32</sub>
        ... (omitted) ...
1226 33
            # for fname in os.listdir('../input/test/'):
1227 34
           for fname in ['1.png']:
1228 35
               test_image = load_image(os.path.join('../input/test', fname))
1229 36
37
               test_x = feature_matrix(test_image)
1230 38
               y_pred, = predict(test_x)
1231 39
               output = y_pred.reshape(test_image.shape) *255.0
1232 40
1233 41
               cv2.imwrite('original_' + fname, test_image)
1234 42 43
               cv2.imwrite('cleaned_' + fname, output)
1235<sub>44</sub>
1236 45
        if __name__ == '__main__':
123746
1238<sup>47</sup>
           main()
1239
```

#### C.1.2 EXAMPLE OF DISCUSSION

1240

1241

# Edge Diffraction in train\_cleaned data

```
1242
        (Lance <TIER: N/A>) I'm studying the pixels in train_cleaned data.&nbsp; I attached a
1243
            colorized blow-up version of part of the image train_cleaned/45.png.   The
            yellow pixels are any pixels that were not pure white ( != 0xFF gray scale) in image 45.
1244
            png, the green was pure white (0xFF).
1245
           So you see what looks like an edge diffraction line lining the outer edge of all the
1246
               letters.
           Okay, maybe I got something wrong in my code.  Can anyone confirm this edge
1247
               diffraction thing in the train_cleaned data, as for example the first word in
               train_cleaned/45.png (There).  You need to make the non-white (byte != 0xFF)
1248
               pixels all a more contrasting color or you may not see it.
1249
           I'm guessing that the clean png files were at some point scanned in using some kind of
               optical scanning machine which added these edge diffraction lines when the light
1250
               diffracts off the edge of the black ink character. 
1251
              (omitted) ...
           + (Rangel Dokov <TIER: MASTER>) Yes, there is some noise, which doesn't look like it
1252
               should be there in the clean set... I ran a test setting everything whiter that 0xF5
1253
               to 0xFF and the RMSE was 0.005, which should be an upper bound on the effects from the
                halos. This will likely be large enough to make the top of the leaderboard a game of
1254
               luck, but since this is just a playground competition I'm not terribly worried about
1255
               it.
1256
```

#### C.2 Example Agent Workflow

1257

1258 1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1276

1277

1278

1279

In our experiment settings, CoMind only accesses top-10 voted discussions and kernels and ignores the rest. The community is initialized with these artifacts. Upon completion of this process, 7 ideas and 10 pipelines are generated. Below is an excerpt of the ideas and reports generated by the Analyzer.

```
(0) Use behaviour-based clustering of neural networks: cluster models by their error patterns
    and ensemble them for document enhancement
(1) Implement sliding-window patch-based models that take an input window and output multiple
    cleaned pixels simultaneously for both denoising and resolution enhancement
```

- (2) Apply a Waifu2x-inspired deep convolutional neural network with gradually increasing filter counts (e.g., 1  $\rightarrow$  32  $\rightarrow$  64  $\rightarrow$  128  $\rightarrow$  256  $\rightarrow$  512  $\rightarrow$  1) and LeakyReLU activations for effective denoising
- (3) Carefully initialize convolutional weights (e.g., stdv = sqrt(2/(kW\*kH\*nOutputPlane))) and use LeakyReLU to improve model convergence and performance
- (4) Ensemble multiple models with different input preprocessing: combine outputs from a pure CNN, background-removed images, edge maps, and thresholded inputs to capture diverse noise characteristics
- (5) Augment training data to simulate real-world 3D deformations and shadows on text, not just 2D noise, to better match test-time artifacts
- (6) Account for systematic artifacts in 'clean' training data (e.g., single-pixel halos) by treating them as noise or adjusting targets accordingly during training

```
Public pipeline (0): - Pipeline: A simple feed-forward neural network that denoises one pixel
    at a time (Theano).
```

- Feature engineering: for each pixel extract a 5\*5 window of gray values (neighbors), 5\*5median blur, 25\*25 median blur, Sobel gradient and second-order derivative magnitudes, stack into a feature vector. Normalize features to [0,1].
- Model architecture: two-layer MLP; hidden layer size N\_HIDDEN=10, tanh activation, output layer with custom activation clip(x+0.5,0,1).
- Training: MSE cost, stochastic gradient descent with learning rate 0.1, batch size 20, 1280 epochs 100. Validation on one image (3.png) at each epoch by RMSE. 1281
  - Prediction: apply same feature\_matrix to test images, predict pixel values, reshape to full image, write out cleaned PNGs.

```
1282
        - Code abstract:
1283
           def feature matrix(img):
              window=(5,5)
              nbrs=[cv2.getRectSubPix(img,window,(y,x)).ravel()
1285
                    for x,y in itertools.product(range(img.shape[0]),range(img.shape[1]))]
              median5=cv2.medianBlur(img,5).ravel()
1286
              median25=cv2.medianBlur(img, 25).ravel()
1287
              grad=np.abs(cv2.Sobel(img,cv2.CV_16S,1,1,ksize=3).ravel())
              div=np.abs(cv2.Sobel(img,cv2.CV_16S,2,2,ksize=3).ravel())
1288
              misc=np.vstack((median5, median25, grad, div)).T
1289
              features=np.hstack((nbrs,misc))
              return (features/255.).astype('float32')
1290
1291
           class Model(object):
              def __init__(...):
1292
                  \verb|self.layer1=Layer(...,n_in=...,n_out=N_HIDDEN,activation=T.tanh)|\\
1293
                  self.layer2=Layer(...,n_in=N_HIDDEN,n_out=n_out,
                                activation=lambda x: T.clip(x+0.5,0,1))
1294
              def cost(self,y): return T.mean((self.output-y)**2)
1295
```

--- PIPELINE SEPARATOR -----

```
Public pipeline (1): - Pipeline: Matching image backgrounds in R (no ML model).
1297
         - Reads test PNGs in batches of 12 images.
         - Flattens each into vectors of size 258*540, stacks as columns.
1298
          - For each pixel location, takes the maximum value across images as an estimate of background
1299
1300
         - Writes out background images as PNG.
        - Code abstract:
1301
           for(i in 1:4) {
            matches=seq(1,205,by=12)+(i-1)*3
1302
             rawData=matrix(0,258*540,length(matches))
1303
             for(j in seg along(matches)){
              imgY=readPNG(file.path(testDir,paste0(matches[j],'.png')))
1304
              rawData[, j] = as.vector(imgY[1:258,1:540])
1305
            background=matrix(apply(rawData,1,max),258,540)
1306
            writePNG(background, paste0('background', matches[j],'.png'))
1307
1308
               --- PIPELINE SEPARATOR --
1309
        Public pipeline (2): - Pipeline: Pixel-wise Random Forest regression (Python, chunk size=1e6).
1310
          - Feature engineering: pad image by mean value (padding=1); extract 3*3 neighborhood per
              pixel, flatten as features.
1311
         - Training data: load all train noisy images, compute features via joblib parallel (n_jobs
              =128), load targets as flattened clean pixel intensities/255.
1312
         - Model: sklearn.ensemble.RandomForestRegressor(warm_start=True, n_jobs=-1). Incrementally
1313
              add one estimator at a time: split training rows into \mathtt{CHUNKSIZE} = 1e6 slices, in each
1314
              slice increase n_{\text{estimators}} by 1 and fit on that slice.
         - Prediction: extract test features similarly, generate index strings "image_row_col",
1315
              predict pixel values, write submission CSV.
        - Code abstract:
1316
           def get_padded(img, padding=1):
1317
             padval=int(round(img.mean()))
               ... return padded
1318
           def get_features_for_image(img,padding=1):
1319
              padded=get_padded(img,padding)
              return np.vstack([padded[i:i+3,j:j+3].reshape(1,-1)
1320
                            for i in range(rows) for j in range(cols)])
1321
           def get_model(X,y):
1322
              model=RandomForestRegressor(n_estimators=0, warm_start=True, n_jobs=-1)
1323
              for start, end in slices:
                 model.set_params(n_estimators=model.n_estimators+1)
1324
                 model.fit(X[start:end],y[start:end])
1325
              return model
1326
        ----- PIPELINE SEPARATOR -----
1327
        ... (omitted) ..
                --- PIPELINE SEPARATOR -----
1328
        Public pipeline (9): - Pipeline: XGBoost with basic features in R.
1329
         - Features: raw pixel value, kmeansThreshold binary, background_Removal feature.
         - Sample 25k pixels, build DMatrix, 5-fold CV, early stop, select rounds, train xgb.
1330
          - Predict on test, clip, build id strings, write submission.csv.
1331
        - Code abstract:
           dat=cbind(y,raw=img2vec(imgX),thresholded=kmeansThreshold(imgX),
1332
                  backgroundRemoval=img2vec(background_Removal(imgX)))
1333
           rows=sample(nrow(dat),25000)
           dtrain=xgb.DMatrix(dat[rows,-1],label=dat[rows,1])
1334
           xgb.tab=xgb.cv(data=dtrain,nfold=5,early.stop.round=50)
1335
           xgb.mod=xgboost(data=dtrain,nrounds=min.err.idx)
           yHat=predict(xgb.mod, newdata=as.matrix(x))
1336
1337
1338
        Then, the Idea Proposer begins brainstorming to propose more original strategies.
1339
        ===SOLUTION PATH 1===
        Multi-Scale Wavelet-U-Net with Attention and Multi-Term Loss
1340
        Description: decompose each image via discrete wavelet transform (DWT) into low-/high-
1341
             frequency bands, process each band in a dual-branch U-Net enhanced with channel- and
1342
             spatial-attention, then invert the wavelet to reconstruct. Train end-to-end with a
             weighted sum of pixel-wise RMSE, wavelet-band consistency, and edge preservation losses.
1343
        - apply 2D DWT to noisy input -> obtain LL, LH, HL, HH sub-bands
        - feed LL into a "coarse" encoder branch, feed concatenated LH/HL/HH into a "detail" encoder
1344
            branch
1345
        - use a U-Net decoder to upsample each branch back to patch size, fuse via learned 1*1
            convolutions
1346
        - insert Convolutional Block Attention Modules (CBAM) after each encoder and decoder block
1347
        - define loss = alpha*pixelRMSE(clean,output) + beta*bandRMSE(wavelet(clean),wavelet(output))
             + gamma*edgeLoss(Sobel(clean),Sobel(output))
1348
        - train on full images with AdamW and a cosine-annealing LR schedule
1349
        ===SOLUTION_PATH_2===
```

```
Stroke-Aware Conditional GAN with OCR-Guided Perceptual Loss
1351
        Description: build a conditional GAN (generator = deep residual encoder-decoder, discriminator
               PatchGAN) that not only minimizes pixel loss but also preserves text strokes-enforce a
1352
              stroke-level loss via a pre-trained small CNN classifier that predicts presence/width of
1353
              strokes. Add an OCR-based perceptual loss: feed predictions through a frozen OCR engine
1354
             embedding and minimize distance to clean embedding.
        - implement generator as ResNet blocks + skip connections (64->128->256->128->64)
1355
        - implement discriminator as 70 \star 70 PatchGAN to focus on local texture
        - include L1 pixel loss + adversarial loss + stroke consistency loss (L1 between stroke-CNN
1356
             features on clean vs. restored)
1357
        - freeze a small text-structure CNN (trained on binary masks) to extract stroke features
        - run Tesseract (or lightweight OCR CNN) on restored vs. clean, extract penultimate-layer
1358
             activations, add perceptual loss term
1359
        - train with R1 gradient penalty and spectral normalization on discriminator
1360
        ===SOLUTION PATH 3===
1361
        Joint Dictionary Learning + Non-Local Patch Aggregation
        Description: learn paired dictionaries (D_noisy, D_clean) for small patches (e.g. 8*8) via
1362
             coupled K-SVD. At test time, extract overlapping patches, compute sparse codes alpha
1363
             under D_noisy via OMP, reconstruct clean patches = D_clean*alpha. Then, perform non-local
              means on the reconstructed patches to exploit self-similarity and average aggregates.
1364
        - sample a large bank of noisy/clean patch pairs, initialize D noisy, D clean with DCT basis
1365
        - run coupled K-SVD to minimize ||D_noisy*alpha - y_noisy|| + ||D_clean*alpha - y_clean|| w.r.
            t. D_noisy, D_clean, alpha
1366
        - at test time, for each image patch y_noisy, compute alpha via Orthogonal Matching Pursuit (
1367
            sparsity <= k)
        - reconstruct y_{clean}_{est} = D_{clean}_{alpha} for each patch
1368
        - perform block-matching to find K nearest patches per reference patch (Euclid dist), stack
1369
             them
        - aggregate reconstructed patches with non-local weights (e.g. Gaussian on reconstruction
1370
             residual)
1371
        ===SOLUTION PATH 4===
1372
        {\tt Self-Supervised\ Blind\ Denoising\ via\ Noise2Void\ +\ Test-Time\ Adaptation}
1373
        Description: exploit purely noisy data-train a small U-Net with masked pixel prediction (
             Noise2Void) on each test image at inference (test-time training). The network learns to
1374
             predict a pixel from its context, gradually adapting to local noise statistics, then you
1375
             run a forward pass to obtain the cleaned image. No clean target needed.
        - define blind-spot or random masking scheme: mask 1% pixels per batch, replace with neighbors
1376
        - build a lightweight CNN (e.g. 5 down/up blocks with skip connections) that predicts a full
1377
             image
        - fine-tune this CNN on each test image for N_iter (e.g. 500 steps) using only masked L2 loss
1378
        - use data augmentation (rotations, flips) on the single test image to diversify contexts
1379
        - after adaptation, perform a clean forward pass without masking to get the denoised output
        - optionally ensemble outputs from multiple random initializations to reduce variance
1380
1381
```

1382

1383

1384 1385

1386

1387

1388

1389

1390

1391

1393

1394

1395

1396

1397

1398

1399

1401

1402

1403

# To remove similar ideas and decompose overlapped ideas, a reconstruction is performed subsequently. 9 ideas are preserved after the filtering and reconstruction. These ideas are then merged with the idea memory.

- (0) Use diverse ensembles by clustering models based on their error patterns and combining outputs from differently preprocessed inputs (e.g., raw, background-removed, edge maps, thresholded) to capture varied noise characteristics.
- (1) Develop a sliding-window, patch-based model that takes an input region and predicts multiple denoised and super-resolved pixels simultaneously for efficient document enhancement.
- (2) Implement a Waifu2x-inspired deep convolutional network with progressively increasing filter counts and LeakyReLU activations, initialized using He normalization for robust convergence in denoising tasks.
- (3) Augment the training set with simulated 3D deformations, lighting variations, and shadows on text to better match real-world scanning artifacts.
  - (4) Identify and compensate for systematic artifacts in the "clean" training targets (e.g., single-pixel halos) by modeling them as noise or adjusting ground-truth labels during training.
  - (5) Build a wavelet-domain dual-branch U-Net: decompose inputs with 2D DWT into approximation and detail bands, encode separately, decode and fuse via learned 1\*1 convolutions, integrate CBAM attention modules, and train end-to-end with a composite loss combining pixel RMSE, wavelet-band RMSE, and edge loss, optimized using AdamW with a cosineannealing LR schedule.
- (6) Apply a GAN-based restoration framework: use a ResNet-based generator and 70\*70 PatchGAN discriminator, train with a combination of L1 pixel loss, adversarial loss, stroke-consistency loss (via a frozen stroke-feature CNN), and perceptual OCR-feature loss, incorporating R1 gradient penalty and spectral normalization.
- (7) Leverage coupled dictionary learning and non-local aggregation: learn noisy and clean dictionaries simultaneously via coupled K-SVD, reconstruct patches using OMP, enhance self-similarity with block-matching, and fuse reconstructed patches with non-local weighting.
- (8) Implement blind-spot test-time adaptation: randomly mask a small percentage of pixels per batch, train a lightweight CNN with skip connections on the single test image using

```
1404
             masked L2 loss and augmentations, fine-tune per image for a few hundred steps, and
1405
             optionally ensemble multiple initializations for variance reduction.
1406
1407
        Based on the ideas received from the Idea Proposer and previous reports, the Coordinator then
1408
        compiles 2 solution drafts.
1409
        Pipeline 1 (Baseline CNN - U-Net with L1 loss and simple augmentations)
1410
        Overview
1411
        - A lightweight U-Net trained on 256 \!\star\!256 patches of noisy/clean document pairs.
        - Single-branch encoder-decoder with skip connections, ReLU activations, batch norm.
1412
        - Loss: L1 pixel-wise loss (less sensitive to outliers than MSE), optional total variation
1413
             regularization.
        - Optimizer: Adam; train on 1 A6000 within 1 hr.
1414
1415
        Data Preparation
        1. Read all noisy (input) and clean (target) train images, normalize intensities to [0,1].
1416
        2. Extract random 256 \star 256 patches (stride = 128) with matching noisy/clean pairs. 3. Data augmentation: random horizontal/vertical flips, +-90 degree rotations.
1417
1418
        4. Create PyTorch DataLoader with batch_size=16 (fits 48 GB) for ^{\sim}100 k patches.
1419
        Network Architecture (PyTorch pseudocode)
1420
        class UNet(nn.Module):
1421
           def __init__(self):
1422
              super().__init__()
               # Encoder
1423
               self.enc1 = DoubleConv(1, 64)
              self.enc2 = Down(64,128)
1424
              self.enc3 = Down(128, 256)
1425
               self.enc4 = Down(256, 512)
               # Bottleneck
1426
              self.bottleneck = Down(512,512)
1427
               # Decoder
              self.up4 = Up(1024, 256)
1428
              self.up3 = Up(512,128)
1429
              self.up2 = Up(256,64)
              self.up1 = Up(128,64)
1430
               self.final = nn.Conv2d(64,1,kernel_size=1)
1431
           def forward(self,x):
              e1=self.enc1(x)
1432
              e2=self.enc2(e1)
1433
              e3=self.enc3(e2)
1434
              e4=self.enc4(e3)
              b = self.bottleneck(e4)
1435
              d4=self.up4(b,e4)
              d3=self.up3(d4,e3)
1436
              d2=self.up2(d3,e2)
1437
              d1=self.up1(d2,e1)
              return torch.sigmoid(self.final(d1))
1438
1439
        Helper modules:
1440
        - DoubleConv = (Conv2d -> BatchNorm2d -> ReLU) *2
        - Down = (MaxPool2d -> DoubleConv)
1441
        - Up = (ConvTranspose2d for upsampling -> concatenate skip -> DoubleConv)
1442
        Training
1443
        - Loss = L1Loss(output, target) + lambda*TV(output) (lambda=1e-5 for smoothness).
        - Optimizer = Adam(lr=1e-3, weight_decay=1e-5).
1444
        - LR schedule: ReduceLROnPlateau(monitor=val_loss, factor=0.5, patience=5).
1445
        - Train for up to 50 epochs; early-stop if val_loss stagnates.
        - Validation: hold out 10% patches to monitor RMSE.
1446
1447
        Inference
        - For each test image (e.g., 540*258), slide 256*256 window with stride=128, predict, and
1448
            average overlapping outputs.
1449
        - Threshold nothing; output raw [0,1] floats per pixel.
1450
        Compute Budget
1451
           100 k patches, batch 16, ~6 k steps per epoch. On A6000: ~2-3 min/epoch => 50 epochs ~ 2 hr
             ; with early stopping < 1 hr.
1452
1453
        Pipeline 2 (Advanced Wavelet U-Net with CBAM and Composite Loss)
1454
        Overview
1455
        - Dual-branch U-Net operating in wavelet domain (Haar DWT) to explicitly denoise tonal and
1456
             textural components.
        - CBAM (Convolutional Block Attention Modules) to adaptively weigh spatial/channel features.
1457
        - Loss = alpha*L1_pixel + beta*L2_wavelet + gamma*EdgeLoss.
        - Optimizer = AdamW + CosineAnnealingLR.
```

```
1459
        Data Preparation
        - Same as Pipeline 1 (256\star256 patches + augmentations).
1460
        - On-the-fly DWT: for each noisy patch, compute one-level Haar DWT -> yields approximation (A)
1461
              and details (H, V, D).
1462
        Network Architecture
1463
        (implemented in PyTorch)
1464
        class WaveletUNet(nn.Module):
1465
           def __init__(self):
1466
             super().__init__()
              # Shared CBAM-Res blocks for Approx and Detail branches
1467
              self.encA1 = CBAMResBlock(1,64)
              self.encD1 = CBAMResBlock(3,64)
1468
              self.pool = nn.MaxPool2d(2)
1469
              self.encA2 = CBAMResBlock(64,128)
              self.encD2 = CBAMResBlock(64,128)
1470
              # Bottleneck
1471
              self.bottleneck = CBAMResBlock (256, 256)
              # Decoder
1472
              self.up2 = UpRes(256,128)
              self.up1 = UpRes(128,64)
1473
              self.final = nn.Conv2d(64,4,kernel_size=1)
1474
           def forward(self, x):
1475
              # x shape: [B,1,H,W]
              A, (H,V,D) = haar_dwt(x) # PyWavelets or custom conv
1476
              d0 = torch.cat([self.encA1(A), self.encD1(torch.cat([H,V,D],dim=1))],dim=1)
1477
              p1 = self.pool(d0)
              {\tt d1 = torch.cat([self.encA2(p1[:, :128]), self.encD2(p1[:, 128:])], dim=1)}
1478
              b = self.bottleneck(self.pool(d1))
1479
              u2 = self.up2(b,d1)
              u1 = self.upl(u2,d0)
1480
              out = self.final(u1) \# [B,4,H,W]
1481
              # Reconstruct denoised x via inverse DWT
              A_{hat}, H_{hat}, V_{hat}, D_{hat} = out[:,0:1], out[:,1:2], out[:,2:3], out[:,3:4]
1482
              return ihaar_dwt((A_hat, (H_hat, V_hat, D_hat)))
1483
        Key building blocks:
1484
        - CBAMResBlock: residual block with Conv-BN-ReLU->Conv-BN, plus CBAM (channel & spatial
             attention after second Conv).
        - UpRes: ConvTranspose2d upsample->concatenate skip->CBAMResBlock.
1486
1487
        - L1_pixel = L1(output, target).
1488
        - L2_wavelet = L2(DWT(output) - DWT(target)).
1489
        - EdgeLoss = L1(Sobel(output), Sobel(target)).
        - TotalLoss = alpha*L1_pixel + beta*L2_wavelet + gamma*EdgeLoss (alpha=1, beta=0.5, gamma=0.2)
1490
1491
       Training
1492
        - Optimizer: AdamW(lr=1e-3, weight_decay=1e-4).
1493
        - Scheduler: CosineAnnealingLR(T_max=50).
        - Batch size=8 (higher memory due to dual branches).
1494
        - Epochs=50, early-stop by composite loss on val set.
1495
        Inference
1496
        - Same sliding-window approach as Pipeline 1.
1497
        Compute Budget
1498
        - CBAM and wavelet layers add ~1.5* compute but within A6000 48 GB: patch batch=8 yields ~250
1499
             steps/epoch, so 50 epochs ~2 hr. Use early-stop or T_max=30 to finish within 1 hr.
1500
1501
        References
        - CBAM: Woo et al., "CBAM: Convolutional Block Attention Module" (ECCV 2018).
1502
        - Haar DWT: PvWavelets or custom conv filters.
1503
        - Edge loss: Sobel via fixed conv kernels.
1504
        Expected Benefits
1505
        - Explicit noise suppression in high-frequency bands.
        - Attention modules focus on text strokes and edge restoration.
1506
        - Composite loss balances global structure and fine details.
1507
```

In our experiment, 2 parallel Coding Agents are invoked to implement the solution drafts respectively. Each instance has a quota of 20 steps, 3 hours in total and 1 hour per execution. An LLM is instructed to monitor the execution process and terminate failed runs.

The best code snippet generated in the first iteration is:

1458

1508

1509

```
1513 1
        import os, glob, random, math
        from PIL import Image
1514 <sup>2</sup>
        import torch
1515 <sup>3</sup>
1516 <sup>4</sup>
        import torch.nn as nn
        import torch.nn.functional as F
1517 <sub>6</sub>
        from torch.utils.data import Dataset, DataLoader
        from torchvision import transforms
1519 8
1520 9
        # -- 1) Haar DWT / inverse DWT --
1521 10 11
        class HaarDWT(nn.Module):
1522 <sub>12</sub>
            def __init__(self):
1523 <sub>13</sub>
                super().__init__()
                wt = torch.zeros(4, 1, 2, 2)
1524 14
                wt[0, 0, :, :] = 0.5
1525 15
                wt[1, 0, 0, 0], wt[1, 0, 0, 1], wt[1, 0, 1, 0], wt[1, 0, 1, 1] = (
1526 16
                   0.5,
1527<sub>18</sub>
                   -0.5,
1528<sub>19</sub>
                   0.5,
                    -0.5,
1529<sup>20</sup>
1530 <sup>21</sup>
1531 <sup>22</sup> <sub>23</sub>
                wt[2, 0, 0, 0], wt[2, 0, 0, 1], wt[2, 0, 1, 0], wt[2, 0, 1, 1] = (
                   0.5,
1532<sub>24</sub>
                   0.5,
1533 25
                   -0.5,
                   -0.5,
1534 26
1535 <sup>27</sup>
                )
                wt[3, 0, 0, 0], wt[3, 0, 0, 1], wt[3, 0, 1, 0], wt[3, 0, 1, 1] = (
1536 28 29
                   0.5,
1537<sub>30</sub>
                    -0.5,
1538 31
                   -0.5,
                   0.5,
1539<sup>32</sup>
                )
1540<sup>33</sup>
1541 34 35
                self.register_buffer("weight", wt)
1542<sub>36</sub>
            def forward(self, x):
1543 37
                return F.conv2d(x, self.weight, stride=2)
1544 38
1545 39
1546 40 41
        class HaarIDWT(nn.Module):
            def __init__(self):
1547<sub>42</sub>
               super().__init__()
1548 43
                wt = HaarDWT().weight.clone()
                self.conv = nn.ConvTranspose2d(4, 1, 2, stride=2, bias=False)
1549 44
1550 45
                self.conv.weight.data.copy_(wt)
                self.conv.weight.requires_grad_(False)
1551 46
47
1552<sub>48</sub>
            def forward(self, coeffs):
1553 49
                return self.conv(coeffs)
1554 50
1555<sup>51</sup>
1556 52
53
        # -- 2) Sobel edge for EdgeLoss --
        class Sobel(nn.Module):
1557 54
            def __init__(self):
1558 55
               super().__init__()
               kx = torch.tensor(
1559 56
                [[1, 0, -1], [2, 0, -2], [1, 0, -1]], dtype=torch.float32).view(1, 1, 3, 3)
1560 <sup>57</sup>
1561 <sup>58</sup> <sub>59</sub>
                ky = torch.tensor(
1562<sub>60</sub>
                    [[1, 2, 1], [0, 0, 0], [-1, -2, -1]], dtype=torch.float32
1563 61
                ).view(1, 1, 3, 3)
                self.register_buffer("wx", kx)
1564 62
1565 63
                self.register_buffer("wy", ky)
    64
```

```
def forward(self, x):
1567<sub>66</sub>
               gx = F.conv2d(x, self.wx, padding=1)
1568<sub>67</sub>
               gy = F.conv2d(x, self.wy, padding=1)
1569 68
               return torch.sqrt(qx * qx + qy * qy + 1e-6)
1570 69
1571 70
        # -- 3) CBAM, ResBlock, UpRes, WaveletUNet --
1572 <sub>72</sub>
       class CBAM(nn.Module):
1573 <sub>73</sub>
           def __init__(self, c, r=16, k=7):
157474
               super().__init__()
               self.mlp = nn.Sequential(
1575<sup>75</sup>
                  nn.Linear(c, c // r, bias=False),
1576 76
                  nn.ReLU(inplace=True),
1577 78
                  nn.Linear(c // r, c, bias=False),
1578<sub>79</sub>
               )
               self.spatial = nn.Conv2d(2, 1, kernel_size=k, padding=k // 2, bias=
1579 80
                   False)
1580
1581 81
1582 <sub>83</sub>
           def forward(self, x):
               b, c, h, w = x.shape
1583<sub>84</sub>
               avg = F.adaptive_avg_pool2d(x, 1).view(b, c)
1584 85
               mx = F.adaptive_max_pool2d(x, 1).view(b, c)
               ca = torch.sigmoid(self.mlp(avg) + self.mlp(mx)).view(b, c, 1, 1)
1585<sup>86</sup>
1586 87
               x2 = x * ca
               avgc = x2.mean(1, True)
1587<sub>89</sub>
               maxc, _ = x2.max(1, True)
1588<sub>90</sub>
               sa = torch.sigmoid(self.spatial(torch.cat([avgc, maxc], 1)))
1589 91
               return x2 * sa
1590<sup>92</sup>
1591 93 94
        class CBAMResBlock(nn.Module):
1592<sub>95</sub>
           def __init__(self, inp, outp):
1593<sub>96</sub>
               super().__init__()
1594 97
               self.conv1 = nn.Conv2d(inp, outp, 3, padding=1, bias=False)
               self.bn1 = nn.BatchNorm2d(outp)
1595 98
1596 99
               self.conv2 = nn.Conv2d(outp, outp, 3, padding=1, bias=False)
1597<sub>101</sub>
               self.bn2 = nn.BatchNorm2d(outp)
               self.cbam = CBAM(outp)
1598102
               self.skip = (
                   nn.Conv2d(inp, outp, 1, bias=False) if inp != outp else nn.
1599103
                       Identity()
1600
1601<sup>104</sup>
               )
               self.relu = nn.ReLU(inplace=True)
1602
1603<sub>107</sub>
           def forward(self, x):
1604108
               idt = self.skip(x)
               out = self.relu(self.bn1(self.conv1(x)))
1605<sup>109</sup>
               out = self.bn2(self.conv2(out))
1606<sup>110</sup>
1607<sub>112</sub>
               out = self.cbam(out) + idt
               return self.relu(out)
160813
160914
        class UpRes(nn.Module):
1610<sup>115</sup>
1611<sup>116</sup>
           def __init__(self, ic, sc, oc):
1612<sub>118</sub>
               super().__init__()
               self.up = nn.ConvTranspose2d(ic, oc, 2, stride=2)
1613119
               self.block = CBAMResBlock(oc + sc, oc)
1614120
           def forward(self, x, skip):
1615<sup>121</sup>
1616<sup>122</sup>
              x = self.up(x)
               if x.shape[-2:] != skip.shape[-2:]:
   123
1617<sub>124</sub>
                   x = F.interpolate(
1618<sub>125</sub>
                      x, size=skip.shape[-2:], mode="bilinear", align_corners=False
1619126
               return self.block(torch.cat([x, skip], 1))
```

```
1620 128
1621<sub>129</sub>
1622<sub>130</sub>
        class WaveletUNet(nn.Module):
            def __init__(self):
162331
1624<sup>132</sup>
                super().__init__()
                self.dwt = HaarDWT()
1625<sup>133</sup>
                self.idwt = HaarIDWT()
1626
                self.eA1 = CBAMResBlock(1, 64)
1627<sub>136</sub>
                self.eD1 = CBAMResBlock(3, 64)
162837
                self.pool = nn.MaxPool2d(2)
                self.eA2 = CBAMResBlock(64, 128)
1629138
                self.eD2 = CBAMResBlock(64, 128)
1630<sup>139</sup>
                self.b = CBAMResBlock(256, 256)
1631<sub>141</sub>
                self.u2 = UpRes(256, 256, 128)
163242
                self.u1 = UpRes(128, 128, 64)
163343
                self.final = nn.Conv2d(64, 4, 1)
1634144
1635<sup>145</sup>
            def forward(self, x):
1635
1636
                A, H, V, D = self.dwt(x).chunk(4, 1)
                a1 = self.eA1(A)
1637148
                d1 = self.eD1(torch.cat([H, V, D], 1))
163849
                d0 = torch.cat([a1, d1], 1)
1639<sup>150</sup>
                p1 = self.pool(d0)
                pA, pD = p1[:, :64], p1[:, 64:]
1640<sup>151</sup>
                a2 = self.eA2(pA)
    152
1641<sub>153</sub>
                d2 = self.eD2(pD)
1642<sub>154</sub>
                d1b = torch.cat([a2, d2], 1)
1643155
                b = self.b(self.pool(d1b))
                u2 = self.u2(b, d1b)
1644<sup>156</sup>
1645<sup>157</sup>
                u1 = self.u1(u2, d0)
                out = self.final(u1)
1646<sub>159</sub>
                return self.idwt(out)
1647<sub>160</sub>
164861
        # -- 4) OCRDataset for full training --
1649<sup>162</sup>
1650<sup>163</sup>
        class OCRDataset(Dataset):
            def __init__(self, noisy, clean, pp, ps, mode):
1651<sub>165</sub>
                self.noisy = noisy
165266
                self.clean = clean
1653167
                self.pp = pp
                self.ps = ps
1654<sup>168</sup>
                self.mode = mode
1655<sup>169</sup>
                self.toT = transforms.ToTensor()
1656<sub>171</sub>
1657<sub>172</sub>
            def __len__(self):
                return len(self.noisy) * self.pp
165873
1659<sup>174</sup>
1660<sup>175</sup>
            def __getitem__(self, idx):
                ii = idx // self.pp
1661<sub>177</sub>
                n = Image.open(self.noisy[ii]).convert("L")
1662 78
                c = Image.open(self.clean[ii]).convert("L")
166379
                tn, tc = self.toT(n), self.toT(c)
1664<sup>180</sup>
                H, W = tn.shape[-2], tn.shape[-1]
                y = random.randint(0, H - self.ps)
x = random.randint(0, W - self.ps)
1665<sup>181</sup>
    182
1666
                tn = tn[:, y : y + self.ps, x : x + self.ps]
1667<sub>184</sub>
                tc = tc[:, y : y + self.ps, x : x + self.ps]
                if self.mode == "train":
1668 85
                    if random.random() < 0.5:</pre>
1669<sup>186</sup>
1670<sup>187</sup>
                        tn, tc = torch.flip(tn, [-1]), torch.flip(tc, [-1])
                    if random.random() < 0.5:</pre>
    188
1671<sub>189</sub>
                       tn, tc = torch.flip(tn, [-2]), torch.flip(tc, [-2])
1672<sub>190</sub>
                    k = random.randint(0, 3)
                    if k > 0:
167391
```

```
1674
1675
                       tn, tc = torch.rot90(tn, k, [1, 2]), torch.rot90(tc, k, [1, 2])
                            21)
1676 93
               return tn, tc
1677194
1678<sup>195</sup>
        # -- 5) Prepare train/val split --
1679<sup>196</sup>
        all_no = sorted(glob.glob("./input/train/*.png"))
1680<sub>198</sub>
        all_cl = sorted(glob.glob("./input/train_cleaned/*.png"))
1681<sub>199</sub>
        idxs = list(range(len(all_no)))
168200
        random.seed(42)
1683201
        random.shuffle(idxs)
1684<sup>202</sup>
        split = int(0.9 * len(idxs))
       tr_idxs, va_idxs = idxs[:split], idxs[split:]
1685,04
       noisy_tr = [all_no[i] for i in tr_idxs]
1686<sub>205</sub>
       clean_tr = [all_cl[i] for i in tr_idxs]
1687206
       noisy_va = [all_no[i] for i in va_idxs]
        clean_va = [all_cl[i] for i in va_idxs]
1688<sup>207</sup>
1689<sup>208</sup>
    209
        # -- 6) DataLoaders --
1690 210
        train_ds = OCRDataset(noisy_tr, clean_tr, pp=25, ps=256, mode="train")
1691<sub>211</sub>
        val_ds = OCRDataset(noisy_va, clean_va, pp=1, ps=256, mode="val")
169212
       train_loader = DataLoader(
1693<sup>213</sup>
            train_ds, batch_size=8, shuffle=True, num_workers=4, pin_memory=True
1694<sup>214</sup>
        val_loader = DataLoader(
    215
1695<sub>216</sub>
           val_ds, batch_size=1, shuffle=False, num_workers=2, pin_memory=True
1696<sub>217</sub>
1697218
        # -- 7) Model, optimizer, scheduler, loss w/ tuned weights --
1698<sup>219</sup>
1699<sup>220</sup>
        device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
       model = WaveletUNet().to(device)
1700,222
        dwt = HaarDWT().to(device)
1701<sub>223</sub>
        sobel = Sobel().to(device)
170224
       opt = torch.optim.AdamW(model.parameters(), lr=1e-3, weight_decay=1e-4)
        sched = torch.optim.lr_scheduler.CosineAnnealingLR(opt, T_max=30)
1703<sup>225</sup>
1704<sup>226</sup>
        # Tuned weights
1705
       alpha = 0.769634
1706<sub>229</sub>
       beta = 0.430968
       gamma = 0.250750
1707230
1708<sup>231</sup>
1709<sup>232</sup>
    233
        def composite_loss(o, t):
1710<sup>233</sup><sub>234</sub>
           11 = F.11_{loss(0, t)}
1711<sub>235</sub>
            12 = F.mse_loss(dwt(o), dwt(t))
            e1 = F.ll_loss(sobel(o), sobel(t))
171236
            return alpha * 11 + beta * 12 + gamma * e1
1713<sup>237</sup>
1714<sup>238</sup>
1715<sub>240</sub>
        # -- 8) Train + validate --
1716<sub>241</sub>
       best_rmse = 1e9
       patience = 5
1717242
1718<sup>243</sup>
        wait = 0
1719<sup>244</sup>
        for epoch in range(1, 31):
    245
           model.train()
1720<sub>246</sub>
            tloss = 0.0
1721<sub>247</sub>
            for xb, yb in train_loader:
               xb, yb = xb.to(device), yb.to(device)
172248
               pred = model(xb)
1723<sup>249</sup>
1724<sup>250</sup>
               loss = composite_loss(pred, yb)
    251
               opt.zero_grad()
1725<sub>252</sub>
               loss.backward()
1726<sub>253</sub>
               opt.step()
               tloss += loss.item()
1727254
            sched.step()
   255
```

```
1728 256
1729<sub>257</sub>
            # sliding-window validation
1730<sub>58</sub>
           model.eval()
1731259
            rmse sum = 0.0
1732<sup>260</sup>
            stride, ps = 128, 256
            with torch.no_grad():
1733<sup>261</sup>
                for nf, cf in zip(noisy_va, clean_va):
    262
1734<sub>263</sub>
                   imn = Image.open(nf).convert("L")
1735<sub>264</sub>
                    imc = Image.open(cf).convert("L")
1736265
                   tn = transforms.ToTensor()(imn).unsqueeze(0).to(device)
1737266
                   tc = transforms.ToTensor()(imc).unsqueeze(0).to(device)
                   _{-}, _{-}, _{-}, _{-} W = tn.shape
1738<sup>267</sup>
                   acc = torch.zeros_like(tn)
1739 269
                   cnt = torch.zeros_like(tn)
1740<sub>270</sub>
                   xs = list(range(0, W - ps + 1, stride)) + [W - ps]
1741271
                   ys = list(range(0, H - ps + 1, stride)) + [H - ps]
                   for x in xs:
1742<sup>272</sup>
1743<sup>273</sup>
                       for y in ys:
    274
                           p = tn[:, :, y : y + ps, x : x + ps]
1744 275
                           out = model(p)
1745<sub>276</sub>
                           acc[:, :, y : y + ps, x : x + ps] += out
1746277
                           cnt[:, :, y : y + ps, x : x + ps] += 1
1747<sup>278</sup>
                   rec = acc / cnt
                   mse = F.mse_loss(rec, tc).item()
1748<sup>279</sup>
                   rmse_sum += math.sqrt(mse)
    280
1749<sub>281</sub>
            val_rmse = rmse_sum / len(noisy_va)
\boldsymbol{1750}_{\!\!282}
            print(
1751283
                f"Epoch {epoch}: train_loss={tloss/len(train_loader):.4f}, val_RMSE
                    ={val_rmse:.6f}"
1752
1753<sup>284</sup>
1754<sub>286</sub>
            if val_rmse < best_rmse:</pre>
1755287
               best_rmse = val_rmse
1756288
               wait = 0
               torch.save(model.state_dict(), "./working1/best.pth")
1757289
1758<sup>290</sup>
    291
                wait += 1
1759 292
               if wait >= patience:
1760293
                   print("Early stopping.")
1761294
                   break
1762<sup>95</sup>
        print("Best validation RMSE:", best_rmse)
1763<sup>296</sup>
1764 298
        # -- 9) Inference & submission --
1765<sub>299</sub>
        model.load_state_dict(torch.load("./working1/best.pth"))
176600
        model.eval()
        test_files = sorted(glob.glob("./input/test/*.png"))
1767<sup>301</sup>
        with open("./working1/submission.csv", "w") as fout:
1768<sup>302</sup>
            fout.write("id, value\n")
1769<sub>304</sub>
            for tf in test_files:
1770<sub>305</sub>
               im = Image.open(tf).convert("L")
1771306
               t = transforms.ToTensor()(im).unsqueeze(0).to(device)
               _, _, H, W = t.shape
1772<sup>307</sup>
1773<sup>308</sup>
               acc = torch.zeros_like(t)
               cnt = torch.zeros_like(t)
1774 310
               xs = list(range(0, W - ps + 1, stride)) + [W - ps]
1775<sub>311</sub>
               ys = list(range(0, H - ps + 1, stride)) + [H - ps]
               with torch.no_grad():
1776312
                   for x in xs:
1777<sup>313</sup>
1778<sup>314</sup>
                       for y in ys:
                           p = t[:, :, y : y + ps, x : x + ps]
1779<sub>316</sub>
                           out = model(p)
1780<sub>317</sub>
                           acc[:, :, y : y + ps, x : x + ps] += out
1781318
                           cnt[:, :, y : y + ps, x : x + ps] += 1
               rec = (acc / cnt).clamp(0, 1).cpu().numpy()[0, 0]
   319
```

```
1782
1783<sup>320</sup>
321
                imq_id = os.path.splitext(os.path.basename(tf))[0]
                for r in range(H):
1784<sub>322</sub>
                    for c in range (W):
178533
                        fout.write(f''\{img\_id\}_{r+1}_{c+1}, \{rec[r,c]:.6f\}\n'')
1786<sup>324</sup>
        print("Submission written to ./working1/submission.csv")
1787<sup>325</sup>
1788
        The Coordinator submits the scripts generated by Coding Agents and their workspace to the com-
1789
        munity. Then, it proceeds to the second iteration. The report generated by the Analyzer for the
1790
        second Coding Agent is:
1791
        Pipeline:
1792
        Our final pipeline is an Advanced Wavelet U-Net with CBAM attention modules, Composite Loss,
1793
             and sliding-window inference. Specifically:
1794
        1. Data Preparation
1795
           - Read all 112 noisy/clean train pairs, split 90/10 train/val by file indices (random seed
1796
               42).
          - On-the-fly extract 25 random 256 \star 256 patches per train image each epoch, with random
1797
               horizontal/vertical flips and 90 degree rotations.
          - Validation uses a single full-image sliding-window pass (256 patches with 50\% overlap) on
1798
               the 10% hold-out.
1799
        2. Model Architecture
1800
           - Haar DWT and inverse DWT implemented as fixed convolutional filters (LL, LH, HL, HH bands)
1801
          - Dual-branch encoder: separate CBAMResBlocks for the approximation (1 channel) and detail
               (3 channels) coefficients.
1803
          - Two-level downsampling with MaxPool, interleaved CBAMResBlocks, and a 256-channel
               bottleneck.
1804
          - Decoder with ConvTranspose upsampling and CBAMResBlock skip-connections (UpRes),
1805
               recombining wavelet bands.
          - Final 1*1 convolution to 4 coefficient channels and inverse DWT back to image domain.
1806
1807
        3. Composite Loss
          - L1(pixels) weighted by alpha=0.769634
1808
          - L2(wavelet coefficients) weighted by beta=0.430968
1809
          - L1(edge Sobel) weighted by gamma=0.250750
1810
        4. Training
1811
           - Optimizer: AdamW(lr=1e-3, weight_decay=1e-4).
          - Scheduler: CosineAnnealingLR(T_max=30).
1812
          - Batch size 8, 30 epochs, early stopping (patience=5) based on full-image RMSE on the val
1813
1814
        5. Inference
1815
           - Sliding-window full-image inference with 256*256 patches, stride=128, averaging overlapped
1816
               pixels.
          - Generate submission.csv with id=image_row_col and float values.
1817
        6. Hyperparameter tuning
1818
          - Optuna over (alpha, beta, gamma) on a 20-train/5-val mini-split for 10 trials, each with 2
1819
               epochs and 10 patches/image.
          - Best found: alpha=0.769634, beta=0.430968, gamma=0.250750.
1820
1821
        Final validation RMSE on 10% hold-out: ~0.012685.
1822
        References: CBAM (Woo et al., ECCV 2018), Haar DWT.
1823
1824
        Summary:
        - Haar DWT:
1825
         Novelty: 2
         Rationale: Uses standard 1-level Haar wavelet decomposition.
1826
         Feasibility: 10
1827
         Rationale: Straightforward fixed conv implementation in PyTorch.
1828
         Effectiveness: 6
         Rationale: Explicit frequency separation aids denoising of high-frequency noise.
1829
         Efficiency: 9
1830
         Rationale: Lightweight fixed kernels on GPU.
         Confidence: 8
1831
         Rationale: Verified on toy and full runs.
        - CBAMResBlock:
1833
         Novelty: 5
1834
         Rationale: Incorporates both channel \& spatial attention for context-aware features.
         Feasibility: 8
1835
          Rationale: Moderate complexity but supported by existing PyTorch modules.
```

Effectiveness: 7

```
1836
         Rationale: Improves focus on text strokes, empirically lowers RMSE.
1837
         Efficiency: 7
         Rationale: Adds overhead (\tilde{1.5}*), but fits within GPU budget.
1838
         Confidence: 9
1839
         Rationale: Confirmed across toy and full training.
1840
        - UpRes (ConvTranspose + CBAM):
1841
         Novelty: 3
         Rationale: Standard upsampling with skip-connections and attention.
1842
         Feasibility: 9
1843
         Rationale: Simple extension of ResBlock.
         Effectiveness: 6
1844
         Rationale: Maintains detail in reconstruction.
1845
         Efficiency: 8
         Rationale: Reasonable GPU use.
1846
         Confidence: 8
1847
         Rationale: Working in end-to-end script.
1848
        - Wavelet UNet overall:
1849
         Novelty: 7
         Rationale: Dual-branch U-Net in wavelet domain combining DWT, CBAM, skip connections.
1850
         Feasibility: 7
1851
         Rationale: Moderate engineering but fully implemented.
         Effectiveness: 8
1852
         Rationale: Achieved RMSE ~0.0127, strong denoising.
1853
         Efficiency: 7
         Rationale: 25 patches*8 batch; 30 epochs in ~20min.
1854
         Confidence: 8
1855
         Rationale: Reproducible on A6000.
1856
        - Composite Loss:
1857
         Novelty: 6
         Rationale: Balances pixel, frequency, and edge domains.
1858
         Feasibility: 8
1859
         Rationale: Simple weighted sum.
         Effectiveness: 8
1860
         Rationale: Lower RMSE than L1 or L2 alone.
1861
         Efficiency: 7
         Rationale: Small overhead for extra transforms.
1862
         Confidence: 9
1863
         Rationale: Verified in tuning and full runs.
1864
        - Data Augmentation:
1865
         Novelty: 3
         Rationale: Flips and rotations.
1866
         Feasibility: 10
1867
         Rationale: Trivial to implement.
         Effectiveness: 5
1868
         Rationale: Standard but helpful for invariance.
1869
         Efficiency: 9
         Rationale: CPU-level only.
1870
         Confidence: 9
1871
         Rationale: Common practice.
1872
        - CosineAnnealingLR & EarlyStopping:
1873
         Novelty: 4
         Rationale: Standard learning-rate decay and val-based stop.
1874
         Feasibility: 10
1875
         Rationale: Built-in PyTorch.
         Effectiveness: 6
1876
         Rationale: Helps convergence and prevents overfit.
1877
         Efficiency: 9
         Rationale: Minimal overhead.
1878
         Confidence: 9
1879
         Rationale: Observed smoothed training curves.
1880
        - Sliding-window Inference:
1881
         Novelty: 3
         Rationale: Ensures full-image coverage.
1882
         Feasibility: 10
1883
         Rationale: Straightforward loops.
1884
         Effectiveness: 7
         Rationale: Accurate RMSE measurement, smooth outputs.
1885
         Efficiency: 8
1886
         Rationale: Reasonable inference time (~10 min).
         Confidence: 9
1887
         Rationale: Verified end-to-end.
1888
        - Hyperparameter Tuning (Optuna):
1889
         Novelty: 4
         Rationale: Automated search of loss weights.
```

```
Feasibility: 8
1891
          Rationale: Low-cost mini-training.
          Effectiveness: 7
          Rationale: Found better weights vs. defaults.
          Efficiency: 5
         Rationale: 10 trials * 2 epochs ~8 min.
          Confidence: 7
1895
         Rationale: Verified chosen parameters.
        Weaknesses and Suggestions:
1897
        1. EMA & TTA Bug: The EMA + 4-way TTA implementation caused RMSE blow-ups, likely due to
             incorrect accumulation. Before reintroducing, validate TTA on toy splits or isolate EMA
             without TTA.
1899
        2. Multi-Level Wavelets: Extend to 2-3 levels of DWT to capture coarser scales, with separate
             branches per level.
1900
        3. Advanced Loss Terms: Incorporate perceptual (VGG) or structural (SSIM) losses to encourage
1901
             visual fidelity.
        4. Dynamic Loss Weights: Schedule (alpha, beta, gamma) during training (e.g., start with high L1
1902
             , gradually emphasize wavelet).
1903
        5. Mixed Precision & Gradient Checkpointing: Further reduce memory and speed up with AMP and
             checkpointing deep layers.
1904
        6. Deeper/Ubiquitous Attention: Explore transformer-style attention (Swin, ViT) in encoder for
1905
              global context.
        7. Ensemble Models: Train multiple random seeds or architectures (ResUNet, UNet++) and average
1906
              outputs.
1907
        8. Multi-Scale Inference: Use patches at multiple scales or whole-image passes via down/
1908
             upscaling for context.
        9. Adaptive Overlap: Use weights at edges of patches to smoothly blend overlaps (Gaussian
1909
             windows).
        10. Distributed Hyperparameter Search: Scale Optuna across more trials or use BOHB for
1910
             architecture search.
1911
        Implementing these will likely push RMSE even lower and improve robustness across document
             noise patterns.
1912
        Full code: ...(omitted)...
1913
```

1914

1915

1916

1917

1918

1919

1920

1921

1922

1923

1924

1925

1926

1927

1928

1929

1930

1931

1932

1933

1934

1935

1936

1937

1938

1939

1940

1941

1942

1943

Following the brainstorming and idea refinement phases, the current idea memory is presented below. Notably, entries (3), (4), (5), and (7) are newly proposed by the agent, while the remaining ideas have been refined from those in the previous iteration.

- (0) Patch-based convolutional super-resolution denoiser: a sliding-window, patch-based model that predicts multiple denoised and super-resolved pixels per patch. Architecture inspired by Waifu2x with progressively increasing filter counts, LeakyReLU activations, and He initialization for robust convergence.
- (1) Advanced wavelet-domain dual-branch U-Net with CBAM attention and composite loss: decompose inputs via fixed 1-level DWT (LL, LH, HL, HH bands), encode approximation and detail separately with CBAM ResBlocks, decode and fuse via 1\*1 convolutions, and train end-to-end using a weighted sum of pixel L1, wavelet-band L2, and edge L1 losses. Optimized with AdamW and cosine-annealing LR scheduling.
- (2) GAN-based restoration framework: a ResNet-based generator and 70\*70 PatchGAN discriminator trained with combined losses-L1 pixel loss, adversarial loss, stroke-consistency loss ( via frozen stroke-feature CNN), and perceptual OCR-feature loss. Includes R1 gradient penalty and spectral normalization for stability.
- (3) Masked autoencoder with vision transformer for denoising: patchify each image into non-overlapping square tokens, randomly mask a high percentage, pretrain a ViT encoder (12 layers, hidden 768, 12 heads) plus light transformer decoder on L2 reconstruction of dirty images, then append an MLP head and fine-tune end-to-end on noisy->clean pairs with L1 pixel + differentiable OCR-confidence loss. Employ random block dropout and color jitter during fine-tuning; at inference use full-image encoding or averaged mask schedules.
- (4) Conditional diffusion-based restoration: define a forward Gaussian-noise diffusion schedule, train a 5-level U-Net conditioned on the dirty image via channel concatenation and FiLM/cross-attention of sinusoidal timestep embeddings. Use the standard DDPM MSE loss with classifier-free guidance, and sample with a deterministic DDIM sampler (\*50 steps). Optionally post-process with bilateral or median filtering to remove speckles.
- (5) Learnable spectral gating in the Fourier domain: compute the 2D FFT of the dirty image, split its spectrum into low/mid/high radial bands, apply learnable complex masks per band, and modulate each by gate scalars predicted by a lightweight CNN on the dirty image. Recombine via inverse FFT and train end-to-end with L2 pixel loss plus a spectral-smoothness regularizer on the masks.
- (6) Hypernetwork-modulated U-Net: extract per-image noise statistics (mean, std, skew, kurtosis, histogram bins), feed into an MLP hypernetwork that outputs FiLM scale (gamma) and shift (beta) parameters for selected convolutional feature maps of a base U-shaped CNN. Randomly augment noise levels during training; train end-to-end on noisy->clean with L1 loss and a small regularizer pushing gamma->1, beta->0. At inference compute stats per image, generate FiLM params, and denoise via the modulated U-Net.
- (7) Blind-spot test-time adaptation: for each test image, randomly mask a subset of pixels and fine-tune a lightweight CNN with skip connections on the single image using masked L2 loss and augmentations for a few hundred gradient steps. Optionally ensemble multiple random initializations to reduce variance.

```
1944
        (8) Multi-model ensemble with diverse preprocessing: cluster trained models by their error
1945
             patterns and combine their outputs. Apply different preprocessing pipelines (raw,
             background-removed, edge maps, thresholded) to the input, denoise with clustered sub-
1946
             ensembles, and fuse predictions for robustness across noise characteristics.
1947
        (9) Enhanced augmentation and target refinement: simulate realistic scanning artifacts by
1948
             applying 3D text deformations, lighting variations, and shadows to clean images. Identify
              and compensate for systematic artifacts in the provided 'clean' targets (e.g., single-
1949
             pixel halos) by either modeling them as noise or adjusting ground-truth labels during
             training.
1950
1951
        And solution drafts generated in this iteration are:
1952
        Pipeline 1: ResNet-34 Encoder U-Net with Multi-Scale Edge & Total-Variation Loss
1953
1954
        Overview:
        A robust baseline using a pretrained ResNet-34 backbone as a U-Net encoder fused with a light-
1955
             weight decoder. Combines L1 loss, Sobel edge loss at multiple scales, and a total-
1956
             variation regularizer to preserve text strokes while smoothing background noise. Mixed
             precision training and sliding-window inference ensure the entire pipeline runs in \tilde{\ }45
1957
             min on an A6000.
1958
        1. Data Preparation
1959
         - Read all train noisy/clean PNGs, normalize to [0,1].
1960
        - Extract on-the-fly 256*256 patches: random crop + random horizontal/vertical flips + 90
             degree rotations.
1961
        - 90/10 split by file indices (seed=42). Use batch size 8-16.
1962
        2. Model Architecture
1963
         - Encoder: torchvision.models.resnet34(pretrained=True), first conv modified to 1->64
1964
             channels.
         - Decoder: four upsampling stages (ConvTranspose2d + Conv2d+BN+ReLU) mirroring ResNet blocks,
1965
              with skip-connections from encoder layers.
        - Final conv 64->1 + Sigmoid.
1966
1967
        3. Loss Function
         Let y_hat and y be predictions and targets.
1968
         - L1Loss(y_hat,y)
1969
         - Edge loss: L1 between Sobel(y_hat) and Sobel(y) at both full resolution and half resolution
1970
              (downsample by 2).
         - TV: lambda*TV(y_hat) where TV = mean(|\nabla xy_hat|+|\nabla yy_hat|).
1971
         Total loss = alpha*L1 + beta*Edge_full + gamma*Edge_half + delta*TV, e.g. alpha=1.0, beta
             =0.5, gamma=0.25, delta=1e-5.
1972
1973
        4. Optimization
         - Optimizer: AdamW(lr=1e-3, weight_decay=1e-4).
1974
         - Scheduler: CosineAnnealingLR(T_max=25).
1975
         - Mixed precision via torch.cuda.amp.
         - Early stopping on validation RMSE (patience=5).
1976
1977
        5. Inference & Submission
        - Perform sliding-window inference on each test image with 256*256 patches, stride=128.
1978
         - Average overlapping patches.
1979
         - Clamp outputs to [0,1], write submission.csv with id=image_row_col.
1980
        Compute budget: ~20 min train + ~5 min inference.
1981
        Pipeline 2: Laplacian-Pyramid Multi-Scale Residual U-Net with Pyramid Loss
1982
1983
        Overview:
        A novel pyramid-domain network that decomposes images into multi-scale Laplacian bands,
1984
             denoises each band via shared-weight residual blocks, and merges them back. Multi-level
1985
             Ll losses focus the model on both coarse structures and fine text details. Efficient and
             fully end-to-end in PyTorch, training finishes in ~50 min on an A6000.
1986
1987
        1. Data & Augmentation
         - Same data split and patch sampling (256*256, flips, rotations).
1988
1989
        2. Laplacian Pyramid Transform (LPT)
         - On the fly, for each patch: create 2-level Gaussian pyramid using avg-pooling (scale
1990
             1->0.5->0.25), then compute Laplacian bands L0=(I-upsample(I/2)), L1=(I/2-upsample(I/4)),
1991
              and a residual low band L2=I/4.
1992
        - Stack [L0,L1,L2] as 3 input channels.
1993
        3. Network
         - Encoder-decoder U-Net with 3-channel input and 3-channel output, 4 down/upsampling stages,
1994
             64->128->256->512->256->128->64 channel widths.
1995
         - Each block: two 3*3 Conv+BN+ReLU and a residual skip connection.
1996
        4. Reconstruction
1997
         - Output is 3 residual corrections R0,R1,R2 to each band.
         - Reconstruct clean patch:
```

```
1998
          a0 = L0 + R0; a1 = L1 + R1; a2 = L2 + R2;
1999
          I = a0 + up(a1) + up(up(a2)).
2000
        5. Loss Function
2001
         - Per-band L1 losses: w0*L1(a0,H0) + w1*L1(a1,H1) + w2*L1(a2,H2), where Hk are ground-truth
             bands and w0=1.0, w1=0.5, w2=0.25.
2002
         - Pixel L1: lambda_pix*L1(I, Igt) with lambda_pix=0.1.
2003
        6. Training
2004
        - AdamW(lr=5e-4, weight_decay=1e-5), one-cycle scheduler (max_lr=5e-4, epochs=30).
2005
         - Batch size 8, mixed precision.
        - Early stopping on val RMSE (patience=6).
2006
2007
        7. Inference
         - On full test images, run sliding-window with 256*256 patches, stride=128.
2008
        - For each patch, apply pyramid -> model -> reconstruct -> accumulate and average.
2009
        - Clamp to [0,1], write submission.csv.
2010
        Compute budget: ~25 min training + ~8 min inference.
2011
2012
        The best code proposed in the second iteration is:
2013
        import os, random, csv
2014 2
        import numpy as np
2015 <sub>3</sub>
        from glob import glob
2016 4
       from PIL import Image
2017 <sup>5</sup>
        import torch
2018 <sup>6</sup>
        import torch.nn as nn
        import torch.nn.functional as F
        from torch.utils.data import Dataset, DataLoader
```

**2044** 37 **2044** 38

**2045** 39

**2046** 40 **2047** 41

**2048** 42 43 **2049** 44

**2050** <sub>45</sub> **2051** <sub>46</sub>

 $T_MAX = 50 \# for LR scheduler$ 

lambda\_mse = 0.1312037280884873 lambda\_ssim = 0.031198904067240532

lambda\_ssim2 = lambda\_ssim / 2

46 # 5) Dataset + augmentations
47 class OCRDataset(Dataset):

# 4) Loss-weight constants (from tuning)

w1, w2, w3, w4 = 1.0, 0.5, 0.25, 1e-5 lambda\_aux = 0.4394633936788146

```
2052 48
           def __init__(self, noisy_list, clean_list, ps, train):
2053 49
              self.noisy, self.clean = noisy_list, clean_list
2054 50
              self.ps, self.train = ps, train
2055 51
              self.to tensor = transforms.ToTensor()
              self.aug = transforms.Compose(
2056 52
2057 53
                      transforms.RandomChoice(
2058 <sub>55</sub>
                         [
2059 56
                             transforms.RandomHorizontalFlip(1.0),
2060 57
                             transforms.RandomVerticalFlip(1.0),
2061 58
                             transforms.RandomRotation(90),
                             transforms.RandomRotation(180),
2062 59
                             transforms.RandomRotation(270),
2063 61
2064 62
                     ),
2065 63
                      transforms.RandomApply([transforms.GaussianBlur(3, (0.1, 2.0)
                          )], p=0.3),
2066
                      transforms.RandomApply([transforms.RandomAdjustSharpness(2.0)
2067 <sup>64</sup>
                          ], p=0.3),
2068 <sub>65</sub>
                  ]
2069<sub>66</sub>
              )
2070 67
2071 <sup>68</sup>
           def __len__(self):
              return len(self.noisy)
2072 69
    70
2073 71
           def __getitem__(self, i):
2074 72
              n = Image.open(self.noisy[i]).convert("L")
2075 73
              c = Image.open(self.clean[i]).convert("L")
              w, h = n.size
2076 74
               # pad
2077 75
              if w < self.ps or h < self.ps:</pre>
2078
                  pad = (0, 0, max(0, self.ps - w), max(0, self.ps - h))
2079<sub>78</sub>
                  n = transforms.functional.pad(n, pad, fill=255)
                  c = transforms.functional.pad(c, pad, fill=255)
2080 79
                  w, h = n.size
2081 80
2082<sup>81</sup>
               # crop
    82
              if self.train:
2083<sub>83</sub>
                  x = random.randint(0, w - self.ps)
2084<sub>84</sub>
                  y = random.randint(0, h - self.ps)
2085 85
                  x = (w - self.ps) // 2
2086 86
2087 87
                  y = (h - self.ps) // 2
              n = n.crop((x, y, x + self.ps, y + self.ps))
2088 89
              c = c.crop((x, y, x + self.ps, y + self.ps))
2089<sub>90</sub>
              if self.train and random.random() < 0.5:</pre>
2090 91
                  n = self.aug(n)
                  c = self.aug(c)
2091 92
              return self.to_tensor(n), self.to_tensor(c)
2092 93
2093<sub>95</sub>
2094 96
       # 6) Prepare train/val split
       noisy_files = sorted(glob(f"{TRAIN_NOISY}/*.png"))
2095 97
       clean_files = [f"{TRAIN_CLEAN}/" + os.path.basename(x) for x in
2096 98
           noisy_files]
2097
       N = len(noisy_files)
2098<sub>100</sub>
       idx = list(range(N))
2099<sub>101</sub>
       random.shuffle(idx)
       ntr = int(0.9 * N)
210002
2101<sup>103</sup>
       tr_idx, va_idx = idx[:ntr], idx[ntr:]
2102<sup>104</sup>
       train_noisy = [noisy_files[i] for i in tr_idx]
       train_clean = [clean_files[i] for i in tr_idx]
   105
2103
       val_noisy = [noisy_files[i] for i in va_idx]
2104<sub>107</sub>
       val_clean = [clean_files[i] for i in va_idx]
210508
       train_ds = OCRDataset(train_noisy, train_clean, PATCH_SIZE, train=True)
```

```
2106
2107<sub>111</sub>
       |val_ds = OCRDataset(val_noisy, val_clean, PATCH_SIZE, train=False)
        train_loader = DataLoader(
2108112
            train_ds, batch_size=BATCH_SIZE, shuffle=True, num_workers=4,
2109
                 pin_memory=True
2110<sup>113</sup>
        val_loader = DataLoader(
2111<sup>114</sup>
            val_ds, batch_size=BATCH_SIZE, shuffle=False, num_workers=4,
2112
                 pin_memory=True
2113<sub>116</sub>
211417
        # 7) Sobel, TV, SSIM helpers
2115118
        sob_x = (
2116<sup>119</sup>
2117 120
            torch.tensor([[1, 0, -1], [2, 0, -2], [1, 0, -1]], dtype=torch.float32
2118121
            .view(1, 1, 3, 3)
            .to(DEVICE)
2119122
2120<sup>123</sup>
2121<sup>124</sup>
        sob_y = sob_x.transpose(2, 3)
2121
2122
126
2123<sub>127</sub>
        def sobel(x):
2124128
            gx = F.conv2d(x, sob_x, padding=1)
            gy = F.conv2d(x, sob_y, padding=1)
2125<sup>129</sup>
            return torch.sqrt(gx * gx + gy * gy + 1e-6)
2126<sup>130</sup>
2127
2128<sub>133</sub>
        def total_variation(x):
212934
            dh = (x[:, :, 1:, :] - x[:, :, :-1, :]).abs().mean()
            dw = (x[:, :, :, 1:] - x[:, :, :, :-1]).abs().mean()
2130<sup>135</sup>
            return dh + dw
2131<sup>136</sup>
2132<sub>38</sub>
2133139
        def ssim_map(a, b, C1=0.01**2, C2=0.03**2):
213440
            mu_a = F.avg_pool2d(a, 3, 1, 1)
            mu_b = F.avg_pool2d(b, 3, 1, 1)
2135<sup>141</sup>
            sa = F.avg_pool2d(a * a, 3, 1, 1) - mu_a * mu_a
sb = F.avg_pool2d(b * b, 3, 1, 1) - mu_b * mu_b
2136<sup>142</sup>
2130
2137
144
            sab = F.avg_pool2d(a * b, 3, 1, 1) - mu_a * mu_b
2138<sub>145</sub>
            num = (2 * mu_a * mu_b + C1) * (2 * sab + C2)
            den = (mu_a * mu_a + mu_b * mu_b + C1) * (sa + sb + C2)
213946
            return num / (den + 1e-8)
2140<sup>147</sup>
2141<sup>148</sup>
2142
        def ssim_loss(a, b):
2143<sub>151</sub>
            return 1.0 - ssim_map(a, b).mean()
214452
2145<sup>153</sup>
        # 8) loss_terms
2146<sup>154</sup>
        11_loss = nn.L1Loss()
2147<sub>156</sub>
        mse_loss = nn.MSELoss()
2148 57
214958
2150<sup>159</sup>
        def loss_terms(pred, target):
2151<sup>160</sup>
            L1v = 11_loss(pred, target)
            MSEv = mse_loss(pred, target)
    161
2152<sub>62</sub>
            Ef = 11_loss(sobel(pred), sobel(target))
2153<sub>63</sub>
            p2, t2 = F.avg_pool2d(pred, 2), F.avg_pool2d(target, 2)
            Eh = 11_{loss(sobel(p2), sobel(t2))}
215464
2155<sup>165</sup>
            TVv = total_variation(pred)
2156<sup>166</sup>
            return L1v, MSEv, Ef, Eh, TVv
2157<sub>168</sub>
2158<sub>169</sub>
        # 9) Model w/ deep supervision
2159170
        class ResUNetDS(nn.Module):
            def ___init___(self):
   171
```

```
2160 172
                super().__init__()
2161<sub>173</sub>
               r34 = tv_models.resnet34(pretrained=True)
2162,74
               self.enc0 = nn.Conv2d(1, 64, 7, 2, 3, bias=False)
2163175
               self.enc0.weight.data = r34.conv1.weight.data.mean(dim=1, keepdim=
2164
                self.bn0, self.relu0, self.pool0 = r34.bn1, r34.relu, r34.maxpool
2165<sup>176</sup>
               self.enc1, self.enc2 = r34.layer1, r34.layer2
2166<sub>178</sub>
               self.enc3, self.enc4 = r34.layer3, r34.layer4
2167<sub>179</sub>
2168 80
               def up(i, o):
2169181
                   return nn.ConvTranspose2d(i, o, 2, 2)
2170<sup>182</sup>
               def cb(i, o):
2171<sub>184</sub>
                   return nn.Sequential(
2172<sub>185</sub>
                       nn.Conv2d(i, o, 3, 1, 1, bias=False),
2173186
                       nn.BatchNorm2d(o),
                       nn.ReLU(inplace=True),
2174<sup>187</sup>
2175<sup>188</sup>
                       nn.Conv2d(o, o, 3, 1, 1, bias=False),
                       nn.BatchNorm2d(o),
2176<sub>190</sub>
                       nn.ReLU(inplace=True),
2177<sub>191</sub>
                   )
217892
2179<sup>193</sup>
                self.up4, self.dec4 = up(512, 256), cb(256 + 256, 256)
                self.up3, self.dec3 = up(256, 128), cb(128 + 128, 128)
2180<sup>194</sup>
    ,
195
                self.up2, self.dec2 = up(128, 64), cb(64 + 64, 64)
2181,196
               self.aux\_up, self.aux\_out = up(64, 64), nn.Conv2d(64, 1, 1)
2182<sub>197</sub>
               self.up1, self.dec1 = up(64, 64), cb(64 + 64, 64)
2183198
               self.up0, self.outc = up(64, 64), nn.Conv2d(64, 1, 1)
               self.sig = nn.Sigmoid()
2184<sup>199</sup>
2185<sup>200</sup>
            def forward(self, x):
2186,202
               x0 = self.relu0(self.bn0(self.enc0(x)))
2187203
               x1 = self.pool0(x0)
2188204
               x2 = self.enc1(x1)
               x3 = self.enc2(x2)
2189<sup>05</sup>
2190<sup>206</sup>
               x4 = self.enc3(x3)
               x5 = self.enc4(x4)
    207
2191 208
219209
               d4 = self.dec4(torch.cat([self.up4(x5), x4], dim=1))
               d3 = self.dec3(torch.cat([self.up3(d4), x3], dim=1))
2193210
2194211
               d2 = self.dec2(torch.cat([self.up2(d3), x2], dim=1))
2195<sup>212</sup>
               aux = self.sig(self.aux_out(self.aux_up(d2)))
               d1 = self.dec1(torch.cat([self.up1(d2), x0], dim=1))
2196<sup>212</sup><sub>214</sub>
               main = self.sig(self.outc(self.up0(d1)))
2197<sub>215</sub>
               return main, aux
219816
2199<sup>217</sup>
2200<sup>218</sup>
        model = ResUNetDS().to(DEVICE)
2201<sub>220</sub>
        # 10) Optimizer, scheduler, scaler
2202<sub>21</sub>
        optimizer = torch.optim.AdamW(model.parameters(), lr=LR, weight_decay=WD)
        scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, T_max=
2203222
             T_MAX)
2204
2205<sup>223</sup>
        scaler = GradScaler()
2206 225
        # 11) Training + snapshot saving
2207<sub>226</sub>
       best_rmse = float("inf")
        patience = 0
220827
        snap\_epochs = set([10, 20, 30, 40, 50])
2209<sup>228</sup>
2210<sup>229</sup>
    230
        for epoch in range(1, MAX_EPOCHS + 1):
2211<sub>231</sub>
           model.train()
2212<sub>32</sub>
            train_loss = 0.0
2213233
            for noisy_img, clean_img in train_loader:
               noisy_img, clean_img = noisy_img.to(DEVICE), clean_img.to(DEVICE)
   234
```

```
2214
235
                optimizer.zero_grad()
2215<sub>236</sub>
                with autocast():
2216<sub>237</sub>
                    main_pred, aux_pred = model(noisy_img)
                    L1v, MSEv, Ef, Eh, TVv = loss_terms(main_pred, clean_img)
2217238
                    s1 = ssim_loss(main_pred, clean_img)
2218<sup>239</sup>
                    p2, t2 = F.avg_pool2d(main_pred, 2), F.avg_pool2d(clean_img, 2)
2219<sup>240</sup>
                    s2 = ssim_loss(p2, t2)
    241
2220 242
                    main_loss = (
2221<sub>243</sub>
                        w1 * L1v
222244
                        + lambda_mse * MSEv
                        + w2 * Ef
2223245
                        + w3 * Eh
2224<sup>246</sup>
                        + w4 * TVv
2225<sub>.48</sub>
                        + lambda_ssim * s1
2226<sub>249</sub>
                        + lambda_ssim2 * s2
2227250
                    )
                    aux_up = F.interpolate(
2228251
2229<sup>252</sup>
                        aux_pred,
    253
                        size=clean_img.shape[-2:],
2230 254
                        mode="bilinear",
2231<sub>255</sub>
                        align_corners=False,
223256
                    La, Ma, Ea, Eh2, TVa = loss_terms(aux_up, clean_img)
2233<sup>257</sup>
                    sa = ssim_loss(aux_up, clean_img)
2234<sup>258</sup>
    259
                    pa, ca = F.avg_pool2d(aux_up, 2), F.avg_pool2d(clean_img, 2)
2235<sub>260</sub>
                    sa2 = ssim_loss(pa, ca)
2236<sub>261</sub>
                    aux_loss = (
2237262
                       w1 * La
                        + lambda_mse * Ma
2238<sup>263</sup>
2239<sup>264</sup>
                        + w2 * Ea
                        + w3 * Eh2
2240<sub>266</sub>
                        + w4 * TVa
2241<sub>267</sub>
                        + lambda_ssim * sa
224268
                        + lambda_ssim2 * sa2
2243<sup>269</sup>
                    )
2244<sup>270</sup>
                    loss = main_loss + lambda_aux * aux_loss
    271
                scaler.scale(loss).backward()
2245<sub>272</sub>
                scaler.step(optimizer)
2246<sub>273</sub>
                scaler.update()
2247274
                train_loss += loss.item()
            scheduler.step()
2248<sup>75</sup>
2249<sup>276</sup>
            # validation
2250 278
            model.eval()
2251<sub>279</sub>
            se, count = 0.0, 0
225280
            with torch.no_grad():
                for noisy_img, clean_img in val_loader:
2253<sup>281</sup>
                    noisy_img, clean_img = noisy_img.to(DEVICE), clean_img.to(DEVICE
2254^{282}
                         )
2255<sub>283</sub>
                    with autocast():
2256284
                       pred, _ = model(noisy_img)
                    se += ((pred - clean_img) ** 2).sum().item()
2257285
2258<sup>286</sup>
                    count += pred.numel()
2259<sup>287</sup>
            val_rmse = np.sqrt(se / count)
            print(
2260 289
                f"Epoch {epoch}: TrainLoss={train_loss/len(train_loader):.4f},
2261
                     ValRMSE={val_rmse:.6f}"
2262290
2263<sup>291</sup>
2264<sup>292</sup>
             # best + snapshot
    293
            if val_rmse < best_rmse:</pre>
2265<sub>294</sub>
                best_rmse = val_rmse
2266<sub>295</sub>
                torch.save(model.state_dict(), os.path.join(WORK_DIR, "best_full.
2267
                     pth"))
                patience = 0
```

```
2268 297
           else:
2269<sub>298</sub>
               patience += 1
2270,99
           if epoch in snap_epochs:
2271300
               torch.save(model.state_dict(), os.path.join(WORK_DIR, f"snap_{epoch
                    }.pth"))
2272
           if patience >= PATIENCE:
2273<sup>301</sup>
               print("Early stopping.")
    302
2274 303
               break
2275<sub>304</sub>
2276305
        print("Best validation RMSE:", best_rmse)
2277806
        # 12) Ensemble load
2278<sup>307</sup>
        ckpts = ["best_full.pth"] + sorted(
2279<sub>309</sub>
           [f for f in os.listdir(WORK_DIR) if f.startswith("snap_")],
2280<sub>310</sub>
           key=lambda x: int(x.split("_")[1].split(".")[0]),
2281311
        ) [-2:]
        ensemble_nets = []
2282312
2283<sup>313</sup>
        for ck in ckpts:
           net = ResUNetDS().to(DEVICE)
2284<sub>315</sub>
           net.load_state_dict(torch.load(os.path.join(WORK_DIR, ck)))
2285<sub>316</sub>
           net.eval()
2286317
           ensemble_nets.append(net)
2287318
2288<sup>319</sup>
        # 13) Sliding-window ensemble inference
    320
2289<sub>321</sub>
        def ensemble_infer(img_arr):
\boldsymbol{2290}_{\!322}
           h, w = img_arr.shape
            inp = torch.from_numpy(img_arr / 255.0).unsqueeze(0).unsqueeze(0).to(
2291323
               DEVICE)
2292
           ph = (PATCH_SIZE - h % STRIDE) % STRIDE
2293<sup>324</sup>
           pw = (PATCH_SIZE - w % STRIDE) % STRIDE
2294326
           inp = F.pad(inp, (0, pw, 0, ph), mode="reflect")
2295327
           _{-}, _{-}, _{-}, _{-} W = inp.shape
2296328
           out = torch.zeros_like(inp)
           wt = torch.zeros_like(inp)
2297<sup>329</sup>
2298<sup>330</sup>
           for y in range(0, H - PATCH_SIZE + 1, STRIDE):
               for x in range(0, W - PATCH_SIZE + 1, STRIDE):
2299<sub>332</sub>
                   patch = inp[:, :, y : y + PATCH_SIZE, x : x + PATCH_SIZE]
2300<sub>333</sub>
                   preds = []
2301334
                   with torch.no_grad(), autocast():
2302<sup>335</sup>
                       for net in ensemble_nets:
                          p, _{-} = net(patch)
2303<sup>336</sup>
                          preds.append(p)
2304<sub>338</sub>
                   avg_p = torch.stack(preds, 0).mean(0)
2305339
                   out[:, :, y : y + PATCH_SIZE, x : x + PATCH_SIZE] += avg_p
                   wt[:, :, y : y + PATCH_SIZE, x : x + PATCH_SIZE] += 1.0
2306340
           out = out / wt
2307841
           out = out[:, :, :h, :w]
2308<sup>342</sup>
           return out.detach().cpu().numpy().squeeze()
2309,44
2310<sub>345</sub>
        # 14) Write submission.csv
231 1346
2312<sup>347</sup>
        submission_path = os.path.join(WORK_DIR, "submission.csv")
2313<sup>348</sup>
        with open (submission_path, "w", newline="") as f:
            writer = csv.writer(f)
2314<sub>350</sub>
           writer.writerow(["id", "value"])
2315,51
           for tf in sorted(
               glob(f"{TEST_DIR}/*.png"), key=lambda x: int(os.path.basename(x).
2316352
                    split(".")[0])
2317
2318<sup>353</sup>
           ):
               img_id = os.path.basename(tf).split(".")[0]
    354
2319<sub>355</sub>
               img = np.array(Image.open(tf).convert("L"), dtype=np.float32)
2320<sub>356</sub>
               den = ensemble_infer(img)
               H, W = den.shape
2321357
               for i in range(H):
   358
```

```
2322 359
                  for j in range (W):
2323
                      writer.writerow([f"{img_id}_{i+1}_{j+1}", f"{den[i,j]:.6f}"])
2324<sub>861</sub>
       print("Submission saved to", submission_path)
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
```