
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, a multi-agent system designed to systematically leverage 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

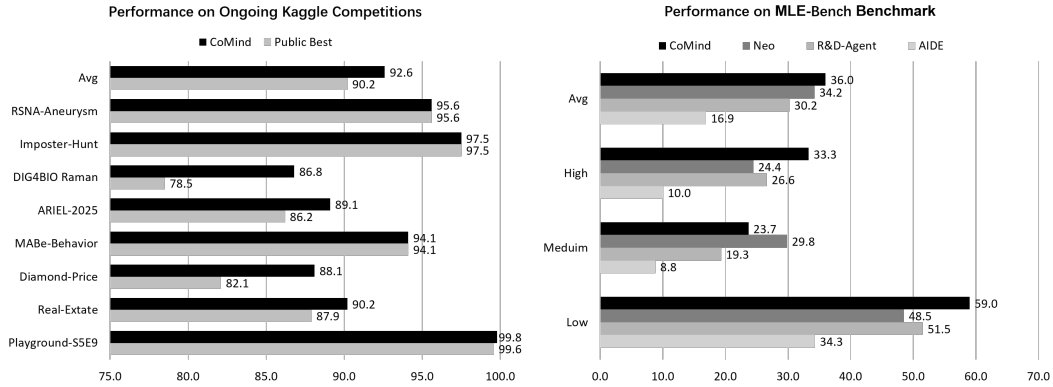


Figure 1: **Left:** CoMind’s win rates on eight ongoing Kaggle competitions compared with the public best entry. **Right:** Any Medal results on 75 MLE-Bench competitions grouped by task difficulty levels. CoMind achieves state-of-the-art performance on MLE-Bench compared to strong baselines.

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* brainstorm 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 utilize 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.

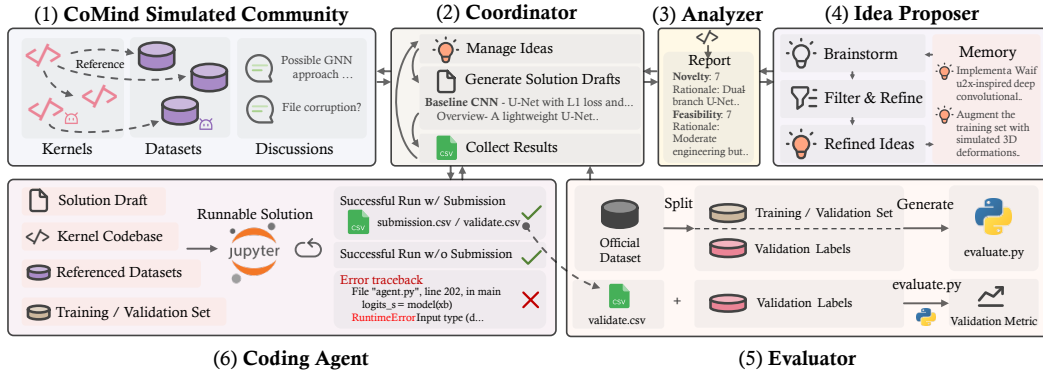


Figure 2: Overview of CoMind. Specialized agents (Coordinator, Analyzer, Idea Proposer, Coding Agent, Evaluator) interact with a simulated Kaggle community of kernels, datasets, and discussions.

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_{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, ReAct-style 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 CoMind 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 leaderboard¹ 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 purpose-built 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 MLEBench, 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-

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

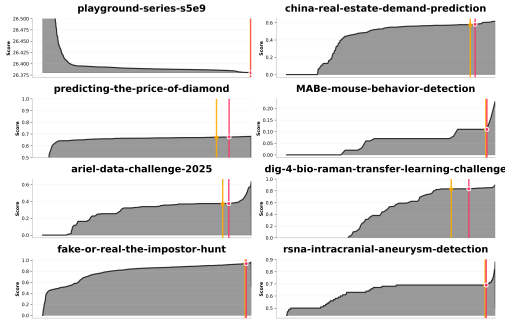


Figure 3: **Left:** Score distributions across participants in eight ongoing Kaggle competitions. Each curve shows the relationship between leaderboard rank (x-axis, inverted) and competition score (y-axis). Vertical lines indicate CoMind’s position (red) and public best performance (yellow). **Right:** 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 I 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.

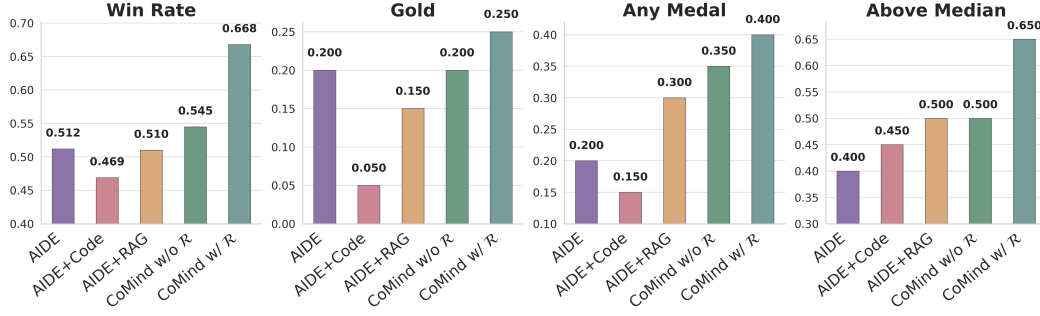


Figure 4: **Performance of CoMind and other baselines on 20 competitions from MLE-Bench-Lite.** *Valid Submission* is the ratio of submissions meeting format requirements and validation criteria. *Win Rate* is the percentage of human competitors outperformed by the agent. *Any Medal*, is the proportion 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 discussions 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.

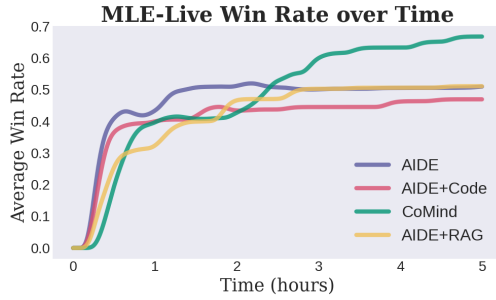


Figure 5: Win rate over time. CoMind sustains improvement while baselines plateau.

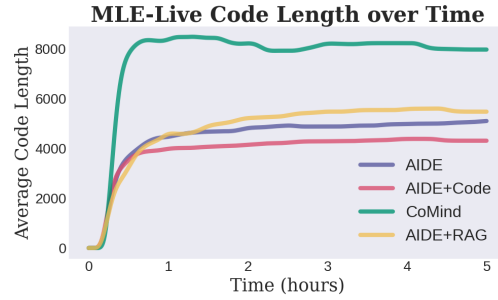


Figure 6: Code complexity over time. CoMind generates longer, richer solutions than baselines.

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 B 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-the-art 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, open-ended 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. MLE-bench: Evaluating machine learning agents on machine learning engineering. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=6s5uXNWGIh>.
- 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. Auto-sklearn 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. Ds-agent: 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 llm 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. MAgentbench: Evaluating language agents on machine learning experimentation. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=1Fs1LvYQW>.

540 Zhengyao Jiang, Dominik Schmidt, Dhruv Srikanth, Dixing Xu, Ian Kaplan, Deniss Jacenko, and
541 Yuxiang Wu. Aide: Ai-driven exploration in the space of code, 2025. URL <https://arxiv.org/abs/2502.13138>.
542
543
544 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
545 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *ArXiv*,
546 abs/2310.06770, 2023a. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:263829697)
547 263829697.
548 Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik
549 Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint*
550 *arXiv:2310.06770*, 2023b.
551
552 Liqiang Jing, Zhehui Huang, Xiaoyang Wang, Wenlin Yao, Wenhao Yu, Kaixin Ma, Hongming
553 Zhang, Xinya Du, and Dong Yu. DS Bench: How far are data science agents from becoming data
554 science experts? In *The Thirteenth International Conference on Learning Representations*, 2025.
555 URL <https://openreview.net/forum?id=DSsSPr0RZJ>.
556
557 Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband:
558 A novel bandit-based approach to hyperparameter optimization. *Journal of Machine Learning*
559 *Research*, 18(185):1–52, 2018. URL <http://jmlr.org/papers/v18/16-558.html>.
560
561 Shanda Li, Tanya Marwah, Junhong Shen, Weiwei Sun, Andrej Risteski, Yiming Yang, and Ameet
562 Talwalkar. Codepde: An inference framework for llm-driven pde solver generation. *arXiv preprint*
563 *arXiv:2505.08783*, 2025.
564
565 Ziming Li, Qianbo Zang, David Ma, Jiawei Guo, Tuney Zheng, Minghao Liu, Xinyao Niu, Yue
566 Wang, Jian Yang, Jiaheng Liu, et al. Autokaggle: A multi-agent framework for autonomous data
567 science competitions. *arXiv preprint arXiv:2410.20424*, 2024.
568
569 Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: Differentiable architecture search. In
570 *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=SleYHoC5FX>.
571
572 Zexi Liu, Yuzhu Cai, Xinyu Zhu, Yujie Zheng, Runkun Chen, Ying Wen, Yanfeng Wang,
573 E Weinan, and Siheng Chen. MI-master: Towards ai-for-ai via integration of exploration and
574 reasoning. *ArXiv*, abs/2506.16499, 2025. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:279465426)
575 *CorpusID:279465426*.
576
577 Peter Mattson, Christine Cheng, Gregory Diamos, Cody Coleman, Paulius Micikevicius, David Pat-
578 terson, Hanlin Tang, Gu-Yeon Wei, Peter Bailis, Victor Bittorf, et al. Mlperf training benchmark.
579 *Proceedings of Machine Learning and Systems*, 2:336–349, 2020.
580
581 OpenAI. Learning to reason with llms, 2024. URL [https://openai.com/index/](https://openai.com/index/learning-to-reason-with-llms/)
582 [learning-to-reason-with-llms/](https://openai.com/index/learning-to-reason-with-llms/).
583
584 OpenAI. Introducing openai o3 and o4-mini, 2025. URL [https://openai.com/index/](https://openai.com/index/introducing-o3-and-o4-mini/)
585 [introducing-o3-and-o4-mini/](https://openai.com/index/introducing-o3-and-o4-mini/).
586
587 Z. Z. Ren, Zhihong Shao, Junxiao Song, Huajian Xin, Haocheng Wang, Wanjia Zhao, Liyue
588 Zhang, Zhe Fu, Qihao Zhu, Dejian Yang, Z. F. Wu, Zhibin Gou, Shirong Ma, Hongxuan Tang,
589 Yuxuan Liu, Wenjun Gao, Daya Guo, and Chong Ruan. Deepseek-prover-v2: Advancing for-
590 mal mathematical reasoning via reinforcement learning for subgoal decomposition, 2025. URL
591 <https://arxiv.org/abs/2504.21801>.
592
593 Bernardino Romera-Paredes, Mohammadamin Barekatin, 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=tN6lDTr4Ed>.
- 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 USE OF LLMs

We employ large language models (LLMs) exclusively for the purpose of assisting in the drafting and refinement of our manuscripts, with the objective of enhancing clarity and coherence.

B 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**. We prioritized these over human annotation to ensure reproducibility and avoid subjective bias.

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
Text Classification	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
	Difficulty	3.76	3.94	1.25	1.31
	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
Audio Classification	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
	Difficulty	7.39	6.46	3.19	6.38
	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
Seq2Seq	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
	Difficulty	3.87	2.14	5.26	3.68
	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
Tabular	CC	2.78	1.62	2.38	0.25
	Pylint Score	8.65	8.96	8.87	9.31
	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

Category	Metric	CoMind	AIDE	AIDE+RAG	AIDE+Code
Image to Image	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
	Difficulty	9.68	6.73	3.94	3.79
	Effort	12 565.66	9967.24	1633.22	1631.93
	SLOC	228.00	175.00	121.00	128.00
	Comment Lines	26.00	8.00	23.00	13.00
	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	11 466.58	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	10 991.00	4841.00	4655.00	5614.00

C TOKEN USAGE ON MLE-BENCH COMPETITIONS

Agent	Uncached Tokens	Cached Tokens	Completion Tokens	Cost
CoMind	18.23 M \pm 14.62 M	13.05 M \pm 11.43 M	1.96 M \pm 1.04 M	\$32.25 \pm 19.43

Table 4: Average token usage (mean \pm standard deviation) and monetary cost for CoMind across the 75 Kaggle competitions in MLE-Bench. The cost is calculated based on the API prices of o4-mini published by OpenAI.

D LOCAL EVALUATION DYNAMICS OF CoMIND ON ONGOING COMPETITIONS

Figure 7 plot the agent’s locally estimated validation metrics over its execution time, as computed by the Evaluator module using held-out validation splits constructed from the public training data. These trajectories capture how CoMind iteratively explores, implements, debugs, and refines solution pipelines through its multi-agent workflow. It illustrates several consistent patterns: (i) rapid early improvements as baseline solutions are assembled; (ii) mid-phase refinements driven by Analyzer-guided idea selection and iterative debugging; and (iii) late-stage stabilization once the agent converges to high-performing configurations.

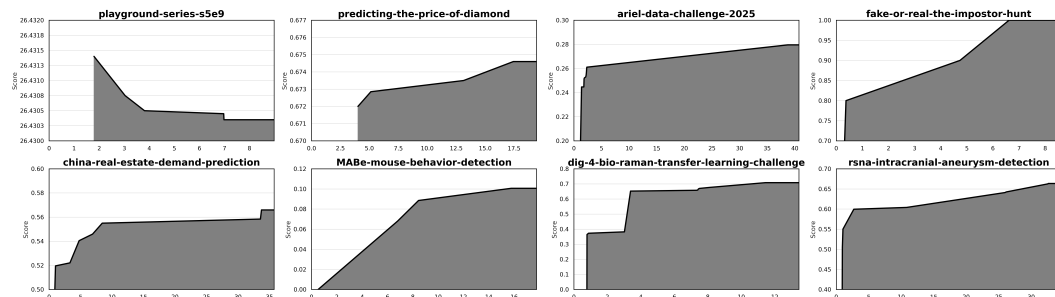


Figure 7: Evolution of CoMind’s Locally Estimated Competition Metrics Over Runtime Across Eight Ongoing Kaggle Tasks.

E CATEGORIES AND DIFFICULTIES IN MLE-BENCH

MLE-Bench (Chan et al., 2025) curates 75 ML engineering-related competitions from Kaggle, creating a diverse set of challenging tasks that test real-world ML engineering skills. These competitions span 15 diverse problem categories. Each competition has an associated description, dataset, and grading code. MLE-Bench categorizes competitions based on human evaluation results: Low (29%) if an experienced ML engineer can produce a sensible solution in 2 hours (excluding training time), Medium (51%) for 2–10 hours, and High (20%) for more than 10 hours.

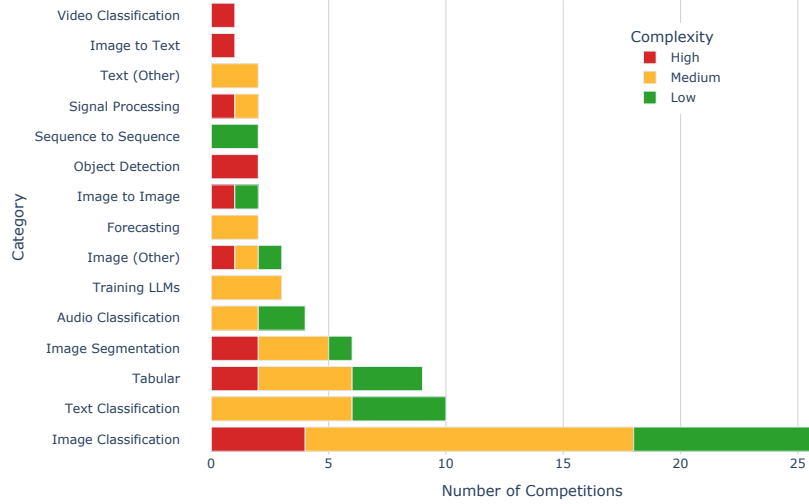


Figure 8: (From MLE-Bench) **Competitions in MLE-Bench spans 15 categories.** MLE-Bench categorizes competition difficulties in Low, Medium and High.

F MLE-LIVE ON MLE-BENCH COMPETITIONS

We collected all public kernels and discussions posted before each competition’s deadline, preserving the temporal information to simulate realistic research environments. Table 5 presents detailed statistics for each competition in MLE-Live. The dataset encompasses 15,733 kernels and 12,951 discussions across competitions of varying difficulty levels (Low, Medium, High).

Figure 9 shows the distribution of kernel votes across all competitions. The distribution is heavily skewed, with over half of the kernels receiving fewer than 10 votes. This long-tail pattern reflects real-world community dynamics where a small fraction of high-quality contributions attracts substantial attention, while most submissions receive minimal engagement.

Table 5: **Statistics of MLE-Live on MLE-Bench competitions.** *Kern.* is the number of public kernels. *Disc.* is the number of discussions. *Dsets* is the number of external datasets referenced by public kernels, *Lines* refers to the average number of lines of all kernels, *Comms* is the average number of comments in all discussions.

Competition	Kern.	Disc.	Dsets	Lines	Comms
Low					
aerial-cactus-identification	275	27	12	177.40	3.52
aptos2019-blindness-detection	186	503	70	357.18	8.33
denoising-dirty-documents	59	19	1	70.36	4.53
detecting-insults-in-social-commentary	0	27	0	—	4.26
dog-breed-identification	64	33	9	161.75	5.15

Continued on next page

Competition	Kern.	Disc.	Dsets	Lines	Comms
dogs-vs-cats-redux-kernels-edition	0	31	0	–	4.68
histopathologic-cancer-detection	118	51	12	228.19	9.53
jigsaw-toxic-comment-classification-challenge	128	394	56	80.88	7.19
leaf-classification	803	23	1	143.37	1.61
mlsp-2013-birds	0	35	0	–	4.94
new-york-city-taxi-fare-prediction	143	53	6	256.69	5.75
nomad2018-predict-transparent-conductors	53	40	9	57.14	5.30
plant-pathology-2020-fgvc7	207	91	27	281.84	5.14
random-acts-of-pizza	8	19	1	35.25	3.74
ranzcr-clip-catheter-line-classification	323	289	156	326.17	6.99
siim-isic-melanoma-classification	276	707	135	345.59	10.87
spooky-author-identification	153	68	13	103.45	3.75
tabular-playground-series-dec-2021	427	134	53	233.22	6.84
tabular-playground-series-may-2022	247	64	23	273.53	5.02
text-normalization-challenge-english-language	57	51	3	153.93	4.45
text-normalization-challenge-russian-language	12	20	2	37.67	2.85
the-icml-2013-whale-challenge-right-whale-redux	0	8	0	–	1.62
Medium					
AI4Code	159	178	93	395.78	5.76
alaska2-image-steganalysis	109	154	31	317.38	6.35
billion-word-imputation	0	23	0	–	3.43
cassava-leaf-disease-classification	411	724	209	358.37	6.77
cdiscoun-image-classification-challenge	91	109	3	57.10	9.25
chaii-hindi-and-tamil-question-answering	269	147	137	341.48	6.74
champs-scalar-coupling	290	340	49	307.36	8.54
facebook-recruiting-iii-keyword-extraction	0	72	0	–	5.32
freesound-audio-tagging-2019	109	128	18	349.46	7.23
google-quest-challenge	225	258	128	442.04	9.04
h-and-m-personalized-fashion-recommendations	419	223	63	240.76	6.30
herbarium-2020-fgvc7	19	15	3	205.83	3.67
herbarium-2021-fgvc8	21	26	17	266.86	2.46
herbarium-2022-fgvc9	36	37	6	274.52	3.84
hotel-id-2021-fgvc8	30	30	15	196.00	2.07
hubmap-kidney-segmentation	321	340	178	325.47	6.70
icecube-neutrinos-in-deep-ice	187	156	38	319.70	5.58
imet-2020-fgvc7	32	21	13	298.37	3.67
inaturalist-2019-fgvc6	12	11	3	244.17	4.27
iwildcam-2020-fgvc7	21	22	5	256.86	3.55
jigsaw-unintended-bias-in-toxicity-classification	410	413	118	301.86	8.16
kuzushiji-recognition	42	55	4	226.57	5.00
learning-agency-lab-automated-essay-scoring-2	477	226	301	349.24	7.78
lmsys-chatbot-arena	305	260	123	299.33	8.51
multi-modal-gesture-recognition	0	39	0	–	3.36
osic-pulmonary-fibrosis-progression	513	499	80	386.23	5.92
petfinder-pawpularity-score	397	461	153	240.50	5.86
plant-pathology-2021-fgvc8	325	112	135	273.29	4.20
seti-breakthrough-listen	230	222	72	333.81	7.82
statoil-iceberg-classifier-challenge	160	178	15	80.98	6.09
tensorflow-speech-recognition-challenge	47	139	4	55.94	6.83
tensorflow2-question-answering	102	193	52	420.29	6.71
tgs-salt-identification-challenge	213	336	24	329.40	12.16
tweet-sentiment-extraction	536	387	117	324.96	9.86
us-patent-phrase-to-phrase-matching	456	233	352	284.25	7.15
uw-madison-gi-tract-image-segmentation	233	230	124	445.30	6.32
ventilator-pressure-prediction	407	268	70	246.07	9.05
whale-categorization-playground	32	22	4	61.93	5.00
High					
3d-object-detection-for-autonomous-vehicles	39	116	10	565.09	6.19
bms-molecular-translation	167	218	57	463.33	10.39

Continued on next page

Competition	Kern.	Disc.	Dsets	Lines	Comms
google-research-identify-contrails-reduce-global-warming	143	136	115	322.43	6.24
hms-harmful-brain-activity-classification	631	390	258	573.88	6.64
iwildcam-2019-fgvc6	37	26	9	196.06	3.77
nfl-player-contact-detection	94	81	32	342.78	4.17
predict-volcanic-eruptions-ingv-oe	102	55	21	271.21	3.55
rsna-2022-cervical-spine-fracture-detection	209	152	106	302.68	4.34
rsna-breast-cancer-detection	538	433	291	342.78	7.11
rsna-miccai-brain-tumor-radiogenomic-classification	414	334	131	321.36	5.72
siim-covid19-detection	586	419	382	337.11	5.59
smartphone-decimeter-2022	49	56	16	244.82	3.79
stanford-covid-vaccine	201	194	29	342.38	8.55
vesuvius-challenge-ink-detection	190	177	72	387.63	6.19
vinbigdata-chest-xray-abnormalities-detection	317	187	115	379.55	6.70

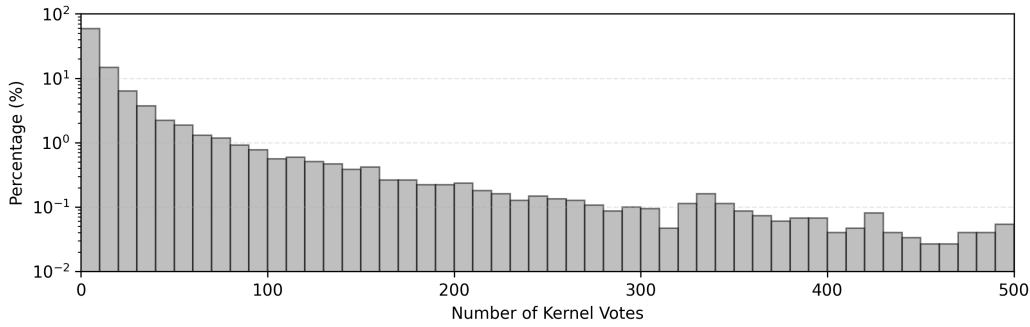


Figure 9: **Distribution of kernel votes across all competitions in MLE-Bench.** Over half of the kernels received fewer than 10 votes, demonstrating the long-tail nature of community contributions where most submissions receive minimal engagement while a small fraction attracts substantial attention.

G ERROR ANALYSIS

In this section, we analyze why CoMind failed to surpass the strongest public kernel in the fake-or-real-the-impostor-hunt Kaggle competition. The task involves identifying the fake text within each pair of samples. Although the dataset originates from The Messenger journal, both the “real” and “fake” texts have been heavily modified by LLMs, making the distinction more subtle and challenging. Below, we summarize the main obstacles encountered during CoMind’s execution:

Noised external resources. Public voting did not reliably indicate kernel quality; many highly upvoted kernels were merely ensembles or reused outputs from stronger solutions, offering little actionable insight. With a large volume of heterogeneous public contributions, identifying genuinely informative resources remained difficult.

Sparse evaluation signal. CoMind depends on the Evaluator’s feedback to guide iteration, but the task’s extremely small dataset allowed validation on only 10 examples. This produced highly unstable feedback and limited the system’s ability to differentiate between small performance variations. As shown in Figure 7, Evaluator accuracy saturated at 100% within the first 7 hours, leaving no meaningful gradient for further improvement.

Insufficient ablation analysis. While CoMind explored multiple strategies, its assessment of individual module effectiveness was inconsistent. For example, early iterations attempted to fine-tune local LLMs via LoRA and train binary classifiers on their outputs, but due to poor base model selection and insufficient hyperparameter exploration, these approaches underperformed simple classical signals such as TF-IDF cosine similarity and token-level similarity.

H 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**.

H.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).
5. 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 ===SEPARATOR=== 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.**

H.2 ANALYZER

Prompt for Strategy Distillation

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 input/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'].

H.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:

1. 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".

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

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 empirical results, not sure whether the evaluation is correct, 10: fully verified on large scale with abundant results)

H.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__.py
  - ...
- 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:

```

```bash
python split_dataset.py --input_dir <path to the input directory> --public_dir ./public
 --private_dir ./private
```

```

DO NOT store the training and test files in other folders such as ./public_`TIMESTAMP`, 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/eval_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. You 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.
```
```

H.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 resources, 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

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

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:

```
<goal>
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>

<code>
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.
</code>

<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}

Consider these factors:

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?

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Respond in the following format:

```
<action>
CONTINUE/STOP
</action>

<explanation>
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.
</explanation>
```

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>

<code>
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.
</code>
```

The validation_submission tag and the submission tag should must be both empty or both non-empty.

I CASE STUDY: DENOISING DIRTY DOCUMENTS

I.1 DATASET PREPARATION

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.

I.1.1 EXAMPLE OF PUBLIC KERNEL

```
1 """
2 A simple feed-forward neural network that denoises one pixel at a time
3 """
4 import numpy as np
5 import theano
6 import theano.tensor as T
7 import cv2
8 import os
9 import itertools
10
11 theano.config.floatX = 'float32'
12
13 def load_image(path):
14     return cv2.imread(path, cv2.IMREAD_GRAYSCALE)
15
16 def feature_matrix(img):
17     """Converts a grayscale image to a feature matrix
18
19     The output value has shape (<number of pixels>, <number of features>)
20     """
21     # select all the pixels in a square around the target pixel as
22     # features
23     window = (5, 5)
24     nbrs = [cv2.getRectSubPix(img, window, (y, x)).ravel()
25             for x, y in itertools.product(range(img.shape[0]), range(img.
26                                     shape[1]))]
27
28     # add some more possibly relevant numbers as features
29     median5 = cv2.medianBlur(img, 5).ravel()
30     median25 = cv2.medianBlur(img, 25).ravel()
31     grad = np.abs(cv2.Sobel(img, cv2.CV_16S, 1, 1, ksize=3).ravel())
32     div = np.abs(cv2.Sobel(img, cv2.CV_16S, 2, 2, ksize=3).ravel())
33
34     ... (omitted) ...
35
36     # for fname in os.listdir('../input/test/'):
37     for fname in ['1.png']:
38         test_image = load_image(os.path.join('../input/test', fname))
39         test_x = feature_matrix(test_image)
40
41         y_pred, = predict(test_x)
42         output = y_pred.reshape(test_image.shape)*255.0
43
44         cv2.imwrite('original_' + fname, test_image)
45         cv2.imwrite('cleaned_' + fname, output)
46
47 if __name__ == '__main__':
48     main()
```

I.1.2 EXAMPLE OF DISCUSSION

```
# Edge Diffraction in train_cleaned data
```



```

1458 Public pipeline (1): - Pipeline: Matching image backgrounds in R (no ML model).
1459   - Reads test PNGs in batches of 12 images.
1460   - Flattens each into vectors of size 258*540, stacks as columns.
1461   - For each pixel location, takes the maximum value across images as an estimate of background
1462   - Writes out background images as PNG.
1463 - Code abstract:
1464   for(i in 1:4) {
1465     matches=seq(1,205,by=12)+(i-1)*3
1466     rawData=matrix(0,258*540,length(matches))
1467     for(j in seq_along(matches)){
1468       imgY=readPNG(file.path(testDir,paste0(matches[j],'.png'))
1469       rawData[,j]=as.vector(imgY[1:258,1:540])
1470     }
1471     background=matrix(apply(rawData,1,max),258,540)
1472     writePNG(background, paste0('background',matches[j],'.png'))
1473   }
1474   ...
1475 ----- PIPELINE SEPARATOR -----
1476 Public pipeline (2): - Pipeline: Pixel-wise Random Forest regression (Python, chunk size=1e6).
1477   - Feature engineering: pad image by mean value (padding=1); extract 3*3 neighborhood per
1478   pixel, flatten as features.
1479   - Training data: load all train noisy images, compute features via joblib parallel (n_jobs
1480   =128), load targets as flattened clean pixel intensities/255.
1481   - Model: sklearn.ensemble.RandomForestRegressor(warm_start=True, n_jobs=-1). Incrementally
1482   add one estimator at a time: split training rows into CHUNKSIZE=1e6 slices, in each
1483   slice increase n_estimators by 1 and fit on that slice.
1484   - Prediction: extract test features similarly, generate index strings "image_row_col",
1485   predict pixel values, write submission CSV.
1486 - Code abstract:
1487   def get_padded(img, padding=1):
1488     padval=int(round(img.mean()))
1489     ... return padded
1490   def get_features_for_image(img,padding=1):
1491     padded=get_padded(img,padding)
1492     return np.vstack([padded[i:i+3,j:j+3].reshape(1,-1)
1493                       for i in range(rows) for j in range(cols)])
1494   ...
1495   def get_model(X,y):
1496     model=RandomForestRegressor(n_estimators=0,warm_start=True,n_jobs=-1)
1497     for start,end in slices:
1498       model.set_params(n_estimators=model.n_estimators+1)
1499       model.fit(X[start:end],y[start:end])
1500     return model
1501   ...
1502 ----- PIPELINE SEPARATOR -----
1503 ... (omitted) ...
1504 ----- PIPELINE SEPARATOR -----
1505 Public pipeline (9): - Pipeline: XGBoost with basic features in R.
1506   - Features: raw pixel value, kmeansThreshold binary, background_Removal feature.
1507   - Sample 25k pixels, build DMatrix, 5-fold CV, early stop, select rounds, train xgb.
1508   - Predict on test, clip, build id strings, write submission.csv.
1509 - Code abstract:
1510   dat=cbind(y,row=img2vec(imgX),thresholded=kmeansThreshold(imgX),
1511             backgroundRemoval=img2vec(background_Removal(imgX)))
1512   rows=sample(nrow(dat),25000)
1513   dtrain=xgb.DMatrix(dat[rows,-1],label=dat[rows,1])
1514   xgb.tab=xgb.cv(data=dtrain,nfold=5,early.stop.round=50)
1515   xgb.mod=xgboost(data=dtrain,nrounds=min.err.idx)
1516   yHat=predict(xgb.mod,newdata=as.matrix(x))
1517   ...
1518
1519 Then, the Idea Proposer begins brainstorming to propose more original strategies.
1520
1521 ===SOLUTION_PATH_1===
1522 Multi-Scale Wavelet-U-Net with Attention and Multi-Term Loss
1523 Description: decompose each image via discrete wavelet transform (DWT) into low-/high-
1524 frequency bands, process each band in a dual-branch U-Net enhanced with channel- and
1525 spatial-attention, then invert the wavelet to reconstruct. Train end-to-end with a
1526 weighted sum of pixel-wise RMSE, wavelet-band consistency, and edge preservation losses.
1527 - apply 2D DWT to noisy input -> obtain LL, LH, HL, HH sub-bands
1528 - feed LL into a "coarse" encoder branch, feed concatenated LH/HL/HH into a "detail" encoder
1529 branch
1530 - use a U-Net decoder to upsample each branch back to patch size, fuse via learned 1*1
1531 convolutions
1532 - insert Convolutional Block Attention Modules (CBAM) after each encoder and decoder block
1533 - define loss = alpha*pixelRMSE(clean,output) + beta*bandRMSE(wavelet(clean),wavelet(output))
1534 + gamma*edgeLoss(Sobel(clean),Sobel(output))
1535 - train on full images with AdamW and a cosine-annealing LR schedule
1536
1537 ===SOLUTION_PATH_2===

```

```

1512 Stroke-Aware Conditional GAN with OCR-Guided Perceptual Loss
1513 Description: build a conditional GAN (generator = deep residual encoder-decoder, discriminator
1514 = PatchGAN) that not only minimizes pixel loss but also preserves text strokes-enforce a
1515 stroke-level loss via a pre-trained small CNN classifier that predicts presence/width of
1516 strokes. Add an OCR-based perceptual loss: feed predictions through a frozen OCR engine
1517 embedding and minimize distance to clean embedding.
1518 - implement generator as ResNet blocks + skip connections (64->128->256->128->64)
1519 - implement discriminator as 70*70 PatchGAN to focus on local texture
1520 - include L1 pixel loss + adversarial loss + stroke consistency loss (L1 between stroke-CNN
1521 features on clean vs. restored)
1522 - freeze a small text-structure CNN (trained on binary masks) to extract stroke features
1523 - run Tesseract (or lightweight OCR CNN) on restored vs. clean, extract penultimate-layer
1524 activations, add perceptual loss term
1525 - train with R1 gradient penalty and spectral normalization on discriminator
1526
1527 ===SOLUTION_PATH_3===
1528 Joint Dictionary Learning + Non-Local Patch Aggregation
1529 Description: learn paired dictionaries (D_noisy, D_clean) for small patches (e.g. 8*8) via
1530 coupled K-SVD. At test time, extract overlapping patches, compute sparse codes alpha
1531 under D_noisy via OMP, reconstruct clean patches = D_clean*alpha. Then, perform non-local
1532 means on the reconstructed patches to exploit self-similarity and average aggregates.
1533 - sample a large bank of noisy/clean patch pairs, initialize D_noisy, D_clean with DCT basis
1534 - run coupled K-SVD to minimize ||D_noisy*alpha - y_noisy|| + ||D_clean*alpha - y_clean|| w.r.
1535 t. D_noisy, D_clean, alpha
1536 - at test time, for each image patch y_noisy, compute alpha via Orthogonal Matching Pursuit (
1537 sparsity <= k)
1538 - reconstruct y_clean_est = D_clean*alpha for each patch
1539 - perform block-matching to find K nearest patches per reference patch (Euclid dist), stack
1540 them
1541 - aggregate reconstructed patches with non-local weights (e.g. Gaussian on reconstruction
1542 residual)
1543
1544 ===SOLUTION_PATH_4===
1545 Self-Supervised Blind Denoising via Noise2Void + Test-Time Adaptation
1546 Description: exploit purely noisy data-train a small U-Net with masked pixel prediction (
1547 Noise2Void) on each test image at inference (test-time training). The network learns to
1548 predict a pixel from its context, gradually adapting to local noise statistics, then you
1549 run a forward pass to obtain the cleaned image. No clean target needed.
1550 - define blind-spot or random masking scheme: mask 1% pixels per batch, replace with neighbors
1551 - build a lightweight CNN (e.g. 5 down/up blocks with skip connections) that predicts a full
1552 image
1553 - fine-tune this CNN on each test image for N_iter (e.g. 500 steps) using only masked L2 loss
1554 - use data augmentation (rotations, flips) on the single test image to diversify contexts
1555 - after adaptation, perform a clean forward pass without masking to get the denoised output
1556 - optionally ensemble outputs from multiple random initializations to reduce variance
1557
1558 To remove similar ideas and decompose overlapped ideas, a reconstruction is performed subse-
1559 quently. 9 ideas are preserved after the filtering and reconstruction. These ideas are then merged
1560 with the idea memory.
1561
1562 (0) Use diverse ensembles by clustering models based on their error patterns and combining
1563 outputs from differently preprocessed inputs (e.g., raw, background-removed, edge maps,
1564 thresholded) to capture varied noise characteristics.
1565 (1) Develop a sliding-window, patch-based model that takes an input region and predicts
1566 multiple denoised and super-resolved pixels simultaneously for efficient document
1567 enhancement.
1568 (2) Implement a Waifu2x-inspired deep convolutional network with progressively increasing
1569 filter counts and LeakyReLU activations, initialized using He normalization for robust
1570 convergence in denoising tasks.
1571 (3) Augment the training set with simulated 3D deformations, lighting variations, and shadows
1572 on text to better match real-world scanning artifacts.
1573 (4) Identify and compensate for systematic artifacts in the "clean" training targets (e.g.,
1574 single-pixel halos) by modeling them as noise or adjusting ground-truth labels during
1575 training.
1576 (5) Build a wavelet-domain dual-branch U-Net: decompose inputs with 2D DWT into approximation
1577 and detail bands, encode separately, decode and fuse via learned 1*1 convolutions,
1578 integrate CBAM attention modules, and train end-to-end with a composite loss combining
1579 pixel RMSE, wavelet-band RMSE, and edge loss, optimized using AdamW with a cosine-
1580 annealing LR schedule.
1581 (6) Apply a GAN-based restoration framework: use a ResNet-based generator and 70*70 PatchGAN
1582 discriminator, train with a combination of L1 pixel loss, adversarial loss, stroke-
1583 consistency loss (via a frozen stroke-feature CNN), and perceptual OCR-feature loss,
1584 incorporating R1 gradient penalty and spectral normalization.
1585 (7) Leverage coupled dictionary learning and non-local aggregation: learn noisy and clean
1586 dictionaries simultaneously via coupled K-SVD, reconstruct patches using OMP, enhance
1587 self-similarity with block-matching, and fuse reconstructed patches with non-local
1588 weighting.
1589 (8) Implement blind-spot test-time adaptation: randomly mask a small percentage of pixels per
1590 batch, train a lightweight CNN with skip connections on the single test image using

```

```

1566         masked L2 loss and augmentations, fine-tune per image for a few hundred steps, and
1567         optionally ensemble multiple initializations for variance reduction.
1568
1569 Based on the ideas received from the Idea Proposer and previous reports, the Coordinator then
1570 compiles 2 solution drafts.
1571
1572 Pipeline 1 (Baseline CNN - U-Net with L1 loss and simple augmentations)
1573 -----
1574 Overview
1575 - A lightweight U-Net trained on 256*256 patches of noisy/clean document pairs.
1576 - Single-branch encoder-decoder with skip connections, ReLU activations, batch norm.
1577 - Loss: L1 pixel-wise loss (less sensitive to outliers than MSE), optional total variation
1578   regularization.
1579 - Optimizer: Adam; train on 1 A6000 within 1 hr.
1580
1581 Data Preparation
1582 1. Read all noisy (input) and clean (target) train images, normalize intensities to [0,1].
1583 2. Extract random 256*256 patches (stride = 128) with matching noisy/clean pairs.
1584 3. Data augmentation: random horizontal/vertical flips, +-90 degree rotations.
1585 4. Create PyTorch DataLoader with batch_size=16 (fits 48 GB) for ~100 k patches.
1586
1587 Network Architecture (PyTorch pseudocode)
1588 ```
1589 class UNet(nn.Module):
1590     def __init__(self):
1591         super().__init__()
1592         # Encoder
1593         self.enc1 = DoubleConv(1, 64)
1594         self.enc2 = Down(64,128)
1595         self.enc3 = Down(128,256)
1596         self.enc4 = Down(256,512)
1597         # Bottleneck
1598         self.bottleneck = Down(512,512)
1599         # Decoder
1600         self.up4 = Up(1024,256)
1601         self.up3 = Up(512,128)
1602         self.up2 = Up(256,64)
1603         self.up1 = Up(128,64)
1604         self.final = nn.Conv2d(64,1,kernel_size=1)
1605     def forward(self,x):
1606         e1=self.enc1(x)
1607         e2=self.enc2(e1)
1608         e3=self.enc3(e2)
1609         e4=self.enc4(e3)
1610         b = self.bottleneck(e4)
1611         d4=self.up4(b,e4)
1612         d3=self.up3(d4,e3)
1613         d2=self.up2(d3,e2)
1614         d1=self.up1(d2,e1)
1615         return torch.sigmoid(self.final(d1))
1616 ```
1617
1618 Helper modules:
1619 - DoubleConv = (Conv2d -> BatchNorm2d -> ReLU) *2
1620 - Down = (MaxPool2d -> DoubleConv)
1621 - Up = (ConvTranspose2d for upsampling -> concatenate skip -> DoubleConv)
1622
1623 Training
1624 - Loss = L1Loss(output, target) + lambda*TV(output) (lambda=1e-5 for smoothness).
1625 - Optimizer = Adam(lr=1e-3, weight_decay=1e-5).
1626 - LR schedule: ReduceLROnPlateau(monitor=val_loss, factor=0.5, patience=5).
1627 - Train for up to 50 epochs; early-stop if val_loss stagnates.
1628 - Validation: hold out 10% patches to monitor RMSE.
1629
1630 Inference
1631 - For each test image (e.g., 540*258), slide 256*256 window with stride=128, predict, and
1632   average overlapping outputs.
1633 - Threshold nothing; output raw [0,1] floats per pixel.
1634
1635 Compute Budget
1636 - ~100 k patches, batch 16, ~6 k steps per epoch. On A6000: ~2-3 min/epoch => 50 epochs ~ 2 hr
1637   ; with early stopping < 1 hr.
1638
1639 Pipeline 2 (Advanced Wavelet U-Net with CBAM and Composite Loss)
1640 -----
1641 Overview
1642 - Dual-branch U-Net operating in wavelet domain (Haar DWT) to explicitly denoise tonal and
1643   textural components.
1644 - CBAM (Convolutional Block Attention Modules) to adaptively weigh spatial/channel features.
1645 - Loss = alpha*L1_pixel + beta*L2_wavelet + gamma*EdgeLoss.
1646 - Optimizer = AdamW + CosineAnnealingLR.

```

```

1620
1621 Data Preparation
1622 - Same as Pipeline 1 (256*256 patches + augmentations).
1623 - On-the-fly DWT: for each noisy patch, compute one-level Haar DWT -> yields approximation (A)
1624   and details (H,V,D).
1625
1626 Network Architecture
1627 (implemented in PyTorch)
1628 ```
1629 class WaveletUNet(nn.Module):
1630     def __init__(self):
1631         super().__init__()
1632         # Shared CBAM-Res blocks for Approx and Detail branches
1633         self.encA1 = CBAMResBlock(1,64)
1634         self.encD1 = CBAMResBlock(3,64)
1635         self.pool = nn.MaxPool2d(2)
1636         self.encA2 = CBAMResBlock(64,128)
1637         self.encD2 = CBAMResBlock(64,128)
1638         # Bottleneck
1639         self.bottleneck = CBAMResBlock(256,256)
1640         # Decoder
1641         self.up2 = UpRes(256,128)
1642         self.up1 = UpRes(128,64)
1643         self.final = nn.Conv2d(64,4,kernel_size=1)
1644     def forward(self, x):
1645         # x shape: [B,1,H,W]
1646         A, (H,V,D) = haar_dwt(x) # PyWavelets or custom conv
1647         d0 = torch.cat([self.encA1(A), self.encD1(torch.cat([H,V,D],dim=1))],dim=1)
1648         p1 = self.pool(d0)
1649         d1 = torch.cat([self.encA2(p1[:, :128]), self.encD2(p1[:, 128:])],dim=1)
1650         b = self.bottleneck(self.pool(d1))
1651         u2 = self.up2(b,d1)
1652         u1 = self.up1(u2,d0)
1653         out = self.final(u1) # [B,4,H,W]
1654         # Reconstruct denoised x via inverse DWT
1655         A_hat, H_hat, V_hat, D_hat = out[:,0:1], out[:,1:2], out[:,2:3], out[:,3:4]
1656         return ihaar_dwt((A_hat, (H_hat,V_hat,D_hat)))
1657
1658 ...
1659
1660 Key building blocks:
1661 - CBAMResBlock: residual block with Conv-BN-ReLU->Conv-BN, plus CBAM (channel & spatial
1662   attention after second Conv).
1663 - UpRes: ConvTranspose2d upsample->concatenate skip->CBAMResBlock.
1664
1665 Loss
1666 - L1_pixel = L1(output, target).
1667 - L2_wavelet = L2(DWT(output) - DWT(target)).
1668 - EdgeLoss = L1(Sobel(output), Sobel(target)).
1669 - TotalLoss = alpha*L1_pixel + beta*L2_wavelet + gamma*EdgeLoss (alpha=1, beta=0.5, gamma=0.2)
1670
1671
1672 Training
1673 - Optimizer: AdamW(lr=1e-3, weight_decay=1e-4).
1674 - Scheduler: CosineAnnealingLR(T_max=50).
1675 - Batch size=8 (higher memory due to dual branches).
1676 - Epochs=50, early-stop by composite loss on val set.
1677
1678 Inference
1679 - Same sliding-window approach as Pipeline 1.
1680
1681 Compute Budget
1682 - CBAM and wavelet layers add ~1.5* compute but within A6000 48 GB: patch batch=8 yields ~250
1683   steps/epoch, so 50 epochs ~2 hr. Use early-stop or T_max=30 to finish within 1 hr.
1684
1685
1686 References
1687 - CBAM: Woo et al., "CBAM: Convolutional Block Attention Module" (ECCV 2018).
1688 - Haar DWT: PyWavelets or custom conv filters.
1689 - Edge loss: Sobel via fixed conv kernels.
1690
1691 Expected Benefits
1692 - Explicit noise suppression in high-frequency bands.
1693 - Attention modules focus on text strokes and edge restoration.
1694 - Composite loss balances global structure and fine details.
1695
1696 In our experiment, 2 parallel Coding Agents are invoked to implement the solution drafts respec-
1697 tively. Each instance has a quota of 20 steps, 3 hours in total and 1 hour per execution. An LLM is
1698 instructed to monitor the execution process and terminate failed runs.
1699
1700 The best code snippet generated in the first iteration is:

```

```

1674
1675 1 import os, glob, random, math
1676 2 from PIL import Image
1677 3 import torch
1678 4 import torch.nn as nn
1679 5 import torch.nn.functional as F
1680 6 from torch.utils.data import Dataset, DataLoader
1681 7 from torchvision import transforms
1682 8
1683 9
1684 10 # -- 1) Haar DWT / inverse DWT --
1685 11 class HaarDWT(nn.Module):
1686 12     def __init__(self):
1687 13         super().__init__()
1688 14         wt = torch.zeros(4, 1, 2, 2)
1689 15         wt[0, 0, :, :] = 0.5
1690 16         wt[1, 0, 0, 0], wt[1, 0, 0, 1], wt[1, 0, 1, 0], wt[1, 0, 1, 1] = (
1691 17             0.5,
1692 18             -0.5,
1693 19             0.5,
1694 20             -0.5,
1695 21         )
1696 22         wt[2, 0, 0, 0], wt[2, 0, 0, 1], wt[2, 0, 1, 0], wt[2, 0, 1, 1] = (
1697 23             0.5,
1698 24             0.5,
1699 25             -0.5,
1700 26             -0.5,
1701 27         )
1702 28         wt[3, 0, 0, 0], wt[3, 0, 0, 1], wt[3, 0, 1, 0], wt[3, 0, 1, 1] = (
1703 29             0.5,
1704 30             -0.5,
1705 31             -0.5,
1706 32             0.5,
1707 33         )
1708 34         self.register_buffer("weight", wt)
1709 35
1710 36     def forward(self, x):
1711 37         return F.conv2d(x, self.weight, stride=2)
1712 38
1713 39
1714 40 class HaarIDWT(nn.Module):
1715 41     def __init__(self):
1716 42         super().__init__()
1717 43         wt = HaarDWT().weight.clone()
1718 44         self.conv = nn.ConvTranspose2d(4, 1, 2, stride=2, bias=False)
1719 45         self.conv.weight.data.copy_(wt)
1720 46         self.conv.weight.requires_grad_(False)
1721 47
1722 48     def forward(self, coeffs):
1723 49         return self.conv(coeffs)
1724 50
1725 51
1726 52 # -- 2) Sobel edge for EdgeLoss --
1727 53 class Sobel(nn.Module):
1728 54     def __init__(self):
1729 55         super().__init__()
1730 56         kx = torch.tensor(
1731 57             [[1, 0, -1], [2, 0, -2], [1, 0, -1]], dtype=torch.float32
1732 58         ).view(1, 1, 3, 3)
1733 59         ky = torch.tensor(
1734 60             [[1, 2, 1], [0, 0, 0], [-1, -2, -1]], dtype=torch.float32
1735 61         ).view(1, 1, 3, 3)
1736 62         self.register_buffer("wx", kx)
1737 63         self.register_buffer("wy", ky)
1738 64

```



```

1728 65 def forward(self, x):
1729 66     gx = F.conv2d(x, self.wx, padding=1)
1730 67     gy = F.conv2d(x, self.wy, padding=1)
1731 68     return torch.sqrt(gx * gx + gy * gy + 1e-6)
1732 69
1733 70 # -- 3) CBAM, ResBlock, UpRes, WaveletUNet --
1734 71
1735 72 class CBAM(nn.Module):
1736 73     def __init__(self, c, r=16, k=7):
1737 74         super().__init__()
1738 75         self.mlp = nn.Sequential(
1739 76             nn.Linear(c, c // r, bias=False),
1740 77             nn.ReLU(inplace=True),
1741 78             nn.Linear(c // r, c, bias=False),
1742 79         )
1743 80         self.spatial = nn.Conv2d(2, 1, kernel_size=k, padding=k // 2, bias=
1744 81             False)
1745 82
1746 83 def forward(self, x):
1747 84     b, c, h, w = x.shape
1748 85     avg = F.adaptive_avg_pool2d(x, 1).view(b, c)
1749 86     mx = F.adaptive_max_pool2d(x, 1).view(b, c)
1750 87     ca = torch.sigmoid(self.mlp(avg) + self.mlp(mx)).view(b, c, 1, 1)
1751 88     x2 = x * ca
1752 89     avgc = x2.mean(1, True)
1753 90     maxc, _ = x2.max(1, True)
1754 91     sa = torch.sigmoid(self.spatial(torch.cat([avgc, maxc], 1)))
1755 92     return x2 * sa
1756 93
1757 94 class CBAMResBlock(nn.Module):
1758 95     def __init__(self, inp, outp):
1759 96         super().__init__()
1760 97         self.conv1 = nn.Conv2d(inp, outp, 3, padding=1, bias=False)
1761 98         self.bn1 = nn.BatchNorm2d(outp)
1762 99         self.conv2 = nn.Conv2d(outp, outp, 3, padding=1, bias=False)
1763 100         self.bn2 = nn.BatchNorm2d(outp)
1764 101         self.cbam = CBAM(outp)
1765 102         self.skip = (
1766 103             nn.Conv2d(inp, outp, 1, bias=False) if inp != outp else nn.
1767 104             Identity()
1768 105         )
1769 106         self.relu = nn.ReLU(inplace=True)
1770 107
1771 108 def forward(self, x):
1772 109     idt = self.skip(x)
1773 110     out = self.relu(self.bn1(self.conv1(x)))
1774 111     out = self.bn2(self.conv2(out))
1775 112     out = self.cbam(out) + idt
1776 113     return self.relu(out)
1777 114
1778 115 class UpRes(nn.Module):
1779 116     def __init__(self, ic, sc, oc):
1780 117         super().__init__()
1781 118         self.up = nn.ConvTranspose2d(ic, oc, 2, stride=2)
1782 119         self.block = CBAMResBlock(oc + sc, oc)
1783 120
1784 121 def forward(self, x, skip):
1785 122     x = self.up(x)
1786 123     if x.shape[-2:] != skip.shape[-2:]:
1787 124         x = F.interpolate(
1788 125             x, size=skip.shape[-2:], mode="bilinear", align_corners=False
1789 126         )
1790 127     return self.block(torch.cat([x, skip], 1))

```

```

1782128
1783129
1784130 class WaveletUNet(nn.Module):
1785131     def __init__(self):
1786132         super().__init__()
1787133         self.dwt = HaarDWT()
1788134         self.idwt = HaarIDWT()
1789135         self.eA1 = CBAMResBlock(1, 64)
1790136         self.eD1 = CBAMResBlock(3, 64)
1791137         self.pool = nn.MaxPool2d(2)
1792138         self.eA2 = CBAMResBlock(64, 128)
1793139         self.eD2 = CBAMResBlock(64, 128)
1794140         self.b = CBAMResBlock(256, 256)
1795141         self.u2 = UpRes(256, 256, 128)
1796142         self.u1 = UpRes(128, 128, 64)
1797143         self.final = nn.Conv2d(64, 4, 1)
1798144
1799145     def forward(self, x):
1800146         A, H, V, D = self.dwt(x).chunk(4, 1)
1801147         a1 = self.eA1(A)
1802148         d1 = self.eD1(torch.cat([H, V, D], 1))
1803149         d0 = torch.cat([a1, d1], 1)
1804150         p1 = self.pool(d0)
1805151         pA, pD = p1[:, :64], p1[:, 64:]
1806152         a2 = self.eA2(pA)
1807153         d2 = self.eD2(pD)
1808154         d1b = torch.cat([a2, d2], 1)
1809155         b = self.b(self.pool(d1b))
1810156         u2 = self.u2(b, d1b)
1811157         u1 = self.u1(u2, d0)
1812158         out = self.final(u1)
1813159         return self.idwt(out)
1814160
1815161
1816162 # -- 4) OCRDataset for full training --
1817163 class OCRDataset(Dataset):
1818164     def __init__(self, noisy, clean, pp, ps, mode):
1819165         self.noisy = noisy
1820166         self.clean = clean
1821167         self.pp = pp
1822168         self.ps = ps
1823169         self.mode = mode
1824170         self.toT = transforms.ToTensor()
1825171
1826172     def __len__(self):
1827173         return len(self.noisy) * self.pp
1828174
1829175     def __getitem__(self, idx):
1830176         ii = idx // self.pp
1831177         n = Image.open(self.noisy[ii]).convert("L")
1832178         c = Image.open(self.clean[ii]).convert("L")
1833179         tn, tc = self.toT(n), self.toT(c)
1834180         H, W = tn.shape[-2], tn.shape[-1]
1835181         y = random.randint(0, H - self.ps)
1836182         x = random.randint(0, W - self.ps)
1837183         tn = tn[:, y : y + self.ps, x : x + self.ps]
1838184         tc = tc[:, y : y + self.ps, x : x + self.ps]
1839185         if self.mode == "train":
1840186             if random.random() < 0.5:
1841187                 tn, tc = torch.flip(tn, [-1]), torch.flip(tc, [-1])
1842188             if random.random() < 0.5:
1843189                 tn, tc = torch.flip(tn, [-2]), torch.flip(tc, [-2])
1844190             k = random.randint(0, 3)
1845191             if k > 0:

```

```

1836         tn, tc = torch.rot90(tn, k, [1, 2]), torch.rot90(tc, k, [1,
1837         2])
1838     return tn, tc
1839
1840 # -- 5) Prepare train/val split --
1841 all_no = sorted(glob.glob("./input/train/*.png"))
1842 all_cl = sorted(glob.glob("./input/train_cleaned/*.png"))
1843 idxs = list(range(len(all_no)))
1844 random.seed(42)
1845 random.shuffle(idxs)
1846 split = int(0.9 * len(idxs))
1847 tr_idx, va_idx = idxs[:split], idxs[split:]
1848 noisy_tr = [all_no[i] for i in tr_idx]
1849 clean_tr = [all_cl[i] for i in tr_idx]
1850 noisy_va = [all_no[i] for i in va_idx]
1851 clean_va = [all_cl[i] for i in va_idx]
1852
1853 # -- 6) DataLoaders --
1854 train_ds = OCRDataset(noisy_tr, clean_tr, pp=25, ps=256, mode="train")
1855 val_ds = OCRDataset(noisy_va, clean_va, pp=1, ps=256, mode="val")
1856 train_loader = DataLoader(
1857     train_ds, batch_size=8, shuffle=True, num_workers=4, pin_memory=True
1858 )
1859 val_loader = DataLoader(
1860     val_ds, batch_size=1, shuffle=False, num_workers=2, pin_memory=True
1861 )
1862
1863 # -- 7) Model, optimizer, scheduler, loss w/ tuned weights --
1864 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
1865 model = WaveletUNet().to(device)
1866 dwt = HaarDWT().to(device)
1867 sobel = Sobel().to(device)
1868 opt = torch.optim.AdamW(model.parameters(), lr=1e-3, weight_decay=1e-4)
1869 sched = torch.optim.lr_scheduler.CosineAnnealingLR(opt, T_max=30)
1870
1871 # Tuned weights
1872 alpha = 0.769634
1873 beta = 0.430968
1874 gamma = 0.250750
1875
1876 def composite_loss(o, t):
1877     l1 = F.l1_loss(o, t)
1878     l2 = F.mse_loss(dwt(o), dwt(t))
1879     e1 = F.l1_loss(sobel(o), sobel(t))
1880     return alpha * l1 + beta * l2 + gamma * e1
1881
1882 # -- 8) Train + validate --
1883 best_rmse = 1e9
1884 patience = 5
1885 wait = 0
1886 for epoch in range(1, 31):
1887     model.train()
1888     tloss = 0.0
1889     for xb, yb in train_loader:
1890         xb, yb = xb.to(device), yb.to(device)
1891         pred = model(xb)
1892         loss = composite_loss(pred, yb)
1893         opt.zero_grad()
1894         loss.backward()
1895         opt.step()
1896         tloss += loss.item()
1897     sched.step()

```

```

1890256
1891257 # sliding-window validation
1892258 model.eval()
1893259 rmse_sum = 0.0
1894260 stride, ps = 128, 256
1895261 with torch.no_grad():
1896262     for nf, cf in zip(noisy_va, clean_va):
1897263         imn = Image.open(nf).convert("L")
1898264         imc = Image.open(cf).convert("L")
1899265         tn = transforms.ToTensor()(imn).unsqueeze(0).to(device)
1900266         tc = transforms.ToTensor()(imc).unsqueeze(0).to(device)
1901267         _, _, H, W = tn.shape
1902268         acc = torch.zeros_like(tn)
1903269         cnt = torch.zeros_like(tn)
1904270         xs = list(range(0, W - ps + 1, stride)) + [W - ps]
1905271         ys = list(range(0, H - ps + 1, stride)) + [H - ps]
1906272         for x in xs:
1907273             for y in ys:
1908274                 p = tn[:, :, y : y + ps, x : x + ps]
1909275                 out = model(p)
1910276                 acc[:, :, y : y + ps, x : x + ps] += out
1911277                 cnt[:, :, y : y + ps, x : x + ps] += 1
1912278             rec = acc / cnt
1913279             mse = F.mse_loss(rec, tc).item()
1914280             rmse_sum += math.sqrt(mse)
1915281 val_rmse = rmse_sum / len(noisy_va)
1916282 print(
1917283     f"Epoch {epoch}: train_loss={tloss/len(train_loader):.4f}, val_RMSE
1918284     ={val_rmse:.6f}"
1919285 )
1920286 if val_rmse < best_rmse:
1921287     best_rmse = val_rmse
1922288     wait = 0
1923289     torch.save(model.state_dict(), "./working1/best.pth")
1924290 else:
1925291     wait += 1
1926292     if wait >= patience:
1927293         print("Early stopping.")
1928294         break
1929295
1930296 print("Best validation RMSE:", best_rmse)
1931297
1932298 # -- 9) Inference & submission --
1933299 model.load_state_dict(torch.load("./working1/best.pth"))
1934300 model.eval()
1935301 test_files = sorted(glob.glob("./input/test/*.png"))
1936302 with open("./working1/submission.csv", "w") as fout:
1937303     fout.write("id,value\n")
1938304     for tf in test_files:
1939305         im = Image.open(tf).convert("L")
1940306         t = transforms.ToTensor()(im).unsqueeze(0).to(device)
1941307         _, _, H, W = t.shape
1942308         acc = torch.zeros_like(t)
1943309         cnt = torch.zeros_like(t)
1944310         xs = list(range(0, W - ps + 1, stride)) + [W - ps]
1945311         ys = list(range(0, H - ps + 1, stride)) + [H - ps]
1946312         with torch.no_grad():
1947313             for x in xs:
1948314                 for y in ys:
1949315                     p = t[:, :, y : y + ps, x : x + ps]
1950316                     out = model(p)
1951317                     acc[:, :, y : y + ps, x : x + ps] += out
1952318                     cnt[:, :, y : y + ps, x : x + ps] += 1
1953319         rec = (acc / cnt).clamp(0, 1).cpu().numpy()[0, 0]

```

```

1944     img_id = os.path.splitext(os.path.basename(tf))[0]
1945320     for r in range(H):
1946321         for c in range(W):
1947322             fout.write(f"{img_id}_{r+1}_{c+1},{rec[r,c]:.6f}\n")
1948323
1949324 print("Submission written to ./working1/submission.csv")

```

The Coordinator submits the scripts generated by Coding Agents and their workspace to the community. Then, it proceeds to the second iteration. The report generated by the Analyzer for the second Coding Agent is:

Pipeline:

Our final pipeline is an Advanced Wavelet U-Net with CBAM attention modules, Composite Loss, and sliding-window inference. Specifically:

1. Data Preparation
 - Read all 112 noisy/clean train pairs, split 90/10 train/val by file indices (random seed 42).
 - On-the-fly extract 25 random 256*256 patches per train image each epoch, with random horizontal/vertical flips and 90 degree rotations.
 - Validation uses a single full-image sliding-window pass (256 patches with 50% overlap) on the 10% hold-out.
2. Model Architecture
 - Haar DWT and inverse DWT implemented as fixed convolutional filters (LL, LH, HL, HH bands).
 - Dual-branch encoder: separate CBAMResBlocks for the approximation (1 channel) and detail (3 channels) coefficients.
 - Two-level downsampling with MaxPool, interleaved CBAMResBlocks, and a 256-channel bottleneck.
 - Decoder with ConvTranspose upsampling and CBAMResBlock skip-connections (UpRes), recombining wavelet bands.
 - Final 1*1 convolution to 4 coefficient channels and inverse DWT back to image domain.
3. Composite Loss
 - L1(pixels) weighted by alpha=0.769634
 - L2(wavelet coefficients) weighted by beta=0.430968
 - L1(edge Sobel) weighted by gamma=0.250750
4. Training
 - Optimizer: AdamW(lr=1e-3, weight_decay=1e-4).
 - Scheduler: CosineAnnealingLR(T_max=30).
 - Batch size 8, 30 epochs, early stopping (patience=5) based on full-image RMSE on the val set.
5. Inference
 - Sliding-window full-image inference with 256*256 patches, stride=128, averaging overlapped pixels.
 - Generate submission.csv with id=image_row_col and float values.
6. Hyperparameter tuning
 - Optuna over (alpha, beta, gamma) on a 20-train/5-val mini-split for 10 trials, each with 2 epochs and 10 patches/image.
 - Best found: alpha=0.769634, beta=0.430968, gamma=0.250750.

Final validation RMSE on 10% hold-out: ~0.012685.

References: CBAM (Woo et al., ECCV 2018), Haar DWT.

Summary:

- Haar DWT:
 - Novelty: 2
 - Rationale: Uses standard 1-level Haar wavelet decomposition.
 - Feasibility: 10
 - Rationale: Straightforward fixed conv implementation in PyTorch.
 - Effectiveness: 6
 - Rationale: Explicit frequency separation aids denoising of high-frequency noise.
 - Efficiency: 9
 - Rationale: Lightweight fixed kernels on GPU.
 - Confidence: 8
 - Rationale: Verified on toy and full runs.
- CBAMResBlock:
 - Novelty: 5
 - Rationale: Incorporates both channel & spatial attention for context-aware features.
 - Feasibility: 8
 - Rationale: Moderate complexity but supported by existing PyTorch modules.
 - Effectiveness: 7

1998 Rationale: Improves focus on text strokes, empirically lowers RMSE.
1999 Efficiency: 7
2000 Rationale: Adds overhead (~1.5*), but fits within GPU budget.
2001 Confidence: 9
2002 Rationale: Confirmed across toy and full training.

2003 - UpRes (ConvTranspose + CBAM):
2004 Novelty: 3
2005 Rationale: Standard upsampling with skip-connections and attention.
2006 Feasibility: 9
2007 Rationale: Simple extension of ResBlock.
2008 Effectiveness: 6
2009 Rationale: Maintains detail in reconstruction.
2010 Efficiency: 8
2011 Rationale: Reasonable GPU use.
2012 Confidence: 8
2013 Rationale: Working in end-to-end script.

2014 - Wavelet UNet overall:
2015 Novelty: 7
2016 Rationale: Dual-branch U-Net in wavelet domain combining DWT, CBAM, skip connections.
2017 Feasibility: 7
2018 Rationale: Moderate engineering but fully implemented.
2019 Effectiveness: 8
2020 Rationale: Achieved RMSE ~0.0127, strong denoising.
2021 Efficiency: 7
2022 Rationale: 25 patches*8 batch; 30 epochs in ~20min.
2023 Confidence: 8
2024 Rationale: Reproducible on A6000.

2025 - Composite Loss:
2026 Novelty: 6
2027 Rationale: Balances pixel, frequency, and edge domains.
2028 Feasibility: 8
2029 Rationale: Simple weighted sum.
2030 Effectiveness: 8
2031 Rationale: Lower RMSE than L1 or L2 alone.
2032 Efficiency: 7
2033 Rationale: Small overhead for extra transforms.
2034 Confidence: 9
2035 Rationale: Verified in tuning and full runs.

2036 - Data Augmentation:
2037 Novelty: 3
2038 Rationale: Flips and rotations.
2039 Feasibility: 10
2040 Rationale: Trivial to implement.
2041 Effectiveness: 5
2042 Rationale: Standard but helpful for invariance.
2043 Efficiency: 9
2044 Rationale: CPU-level only.
2045 Confidence: 9
2046 Rationale: Common practice.

2047 - CosineAnnealingLR & EarlyStopping:
2048 Novelty: 4
2049 Rationale: Standard learning-rate decay and val-based stop.
2050 Feasibility: 10
2051 Rationale: Built-in PyTorch.
Effectiveness: 6
Rationale: Helps convergence and prevents overfit.
Efficiency: 9
Rationale: Minimal overhead.
Confidence: 9
Rationale: Observed smoothed training curves.

2052 - Sliding-window Inference:
2053 Novelty: 3
2054 Rationale: Ensures full-image coverage.
2055 Feasibility: 10
2056 Rationale: Straightforward loops.
2057 Effectiveness: 7
2058 Rationale: Accurate RMSE measurement, smooth outputs.
2059 Efficiency: 8
2060 Rationale: Reasonable inference time (~10 min).
2061 Confidence: 9
2062 Rationale: Verified end-to-end.

2063 - Hyperparameter Tuning (Optuna):
2064 Novelty: 4
2065 Rationale: Automated search of loss weights.

2052 Feasibility: 8
2053 Rationale: Low-cost mini-training.
2054 Effectiveness: 7
2055 Rationale: Found better weights vs. defaults.
2056 Efficiency: 5
2057 Rationale: 10 trials * 2 epochs ~8 min.
2058 Confidence: 7
2059 Rationale: Verified chosen parameters.

2058 Weaknesses and Suggestions:

2059 1. EMA & TTA Bug: The EMA + 4-way TTA implementation caused RMSE blow-ups, likely due to
2060 incorrect accumulation. Before reintroducing, validate TTA on toy splits or isolate EMA
2061 without TTA.

2062 2. Multi-Level Wavelets: Extend to 2-3 levels of DWT to capture coarser scales, with separate
2063 branches per level.

2064 3. Advanced Loss Terms: Incorporate perceptual (VGG) or structural (SSIM) losses to encourage
2065 visual fidelity.

2066 4. Dynamic Loss Weights: Schedule (alpha,beta,gamma) during training (e.g., start with high L1
2067 , gradually emphasize wavelet).

2068 5. Mixed Precision & Gradient Checkpointing: Further reduce memory and speed up with AMP and
2069 checkpointing deep layers.

2070 6. Deeper/Ubiquitous Attention: Explore transformer-style attention (Swin, ViT) in encoder for
2071 global context.

2072 7. Ensemble Models: Train multiple random seeds or architectures (ResUNet, UNet++) and average
2073 outputs.

2074 8. Multi-Scale Inference: Use patches at multiple scales or whole-image passes via down/
2075 upscaling for context.

2076 9. Adaptive Overlap: Use weights at edges of patches to smoothly blend overlaps (Gaussian
2077 windows).

2078 10. Distributed Hyperparameter Search: Scale Optuna across more trials or use BOHB for
2079 architecture search.

2080 Implementing these will likely push RMSE even lower and improve robustness across document
2081 noise patterns.

2082 Full code: ...(omitted)...

2083

2084 Following the brainstorming and idea refinement phases, the current idea memory is presented be-
2085 low. Notably, entries (3), (4), (5), and (7) are newly proposed by the agent, while the remaining
2086 ideas have been refined from those in the previous iteration.

2087

2088 (0) Patch-based convolutional super-resolution denoiser: a sliding-window, patch-based model
2089 that predicts multiple denoised and super-resolved pixels per patch. Architecture
2090 inspired by Waifu2x with progressively increasing filter counts, LeakyReLU activations,
2091 and He initialization for robust convergence.

2092 (1) Advanced wavelet-domain dual-branch U-Net with CBAM attention and composite loss:
2093 decompose inputs via fixed 1-level DWT (LL, LH, HL, HH bands), encode approximation and
2094 detail separately with CBAM ResBlocks, decode and fuse via 1*1 convolutions, and train
2095 end-to-end using a weighted sum of pixel L1, wavelet-band L2, and edge L1 losses.
2096 Optimized with AdamW and cosine-annealing LR scheduling.

2097 (2) GAN-based restoration framework: a ResNet-based generator and 70*70 PatchGAN discriminator
2098 trained with combined losses-L1 pixel loss, adversarial loss, stroke-consistency loss (
2099 via frozen stroke-feature CNN), and perceptual OCR-feature loss. Includes R1 gradient
2100 penalty and spectral normalization for stability.

2101 (3) Masked autoencoder with vision transformer for denoising: patchify each image into non-
2102 overlapping square tokens, randomly mask a high percentage, pretrain a ViT encoder (12
2103 layers, hidden 768, 12 heads) plus light transformer decoder on L2 reconstruction of
2104 dirty images, then append an MLP head and fine-tune end-to-end on noisy->clean pairs with
2105 L1 pixel + differentiable OCR-confidence loss. Employ random block dropout and color
2106 jitter during fine-tuning; at inference use full-image encoding or averaged mask
2107 schedules.

2108 (4) Conditional diffusion-based restoration: define a forward Gaussian-noise diffusion
2109 schedule, train a 5-level U-Net conditioned on the dirty image via channel concatenation
2110 and FiLM/cross-attention of sinusoidal timestep embeddings. Use the standard DDPM MSE
2111 loss with classifier-free guidance, and sample with a deterministic DDIM sampler (~50
2112 steps). Optionally post-process with bilateral or median filtering to remove speckles.

2113 (5) Learnable spectral gating in the Fourier domain: compute the 2D FFT of the dirty image,
2114 split its spectrum into low/mid/high radial bands, apply learnable complex masks per band
2115 , and modulate each by gate scalars predicted by a lightweight CNN on the dirty image.
2116 Recombine via inverse FFT and train end-to-end with L2 pixel loss plus a spectral-
2117 smoothness regularizer on the masks.

2118 (6) Hypernetwork-modulated U-Net: extract per-image noise statistics (mean, std, skew,
2119 kurtosis, histogram bins), feed into an MLP hypernetwork that outputs FiLM scale (gamma)
2120 and shift (beta) parameters for selected convolutional feature maps of a base U-shaped
2121 CNN. Randomly augment noise levels during training; train end-to-end on noisy->clean with
2122 L1 loss and a small regularizer pushing gamma->1, beta->0. At inference compute stats
2123 per image, generate FiLM params, and denoise via the modulated U-Net.

2124 (7) Blind-spot test-time adaptation: for each test image, randomly mask a subset of pixels and
2125 fine-tune a lightweight CNN with skip connections on the single image using masked L2
2126 loss and augmentations for a few hundred gradient steps. Optionally ensemble multiple
2127 random initializations to reduce variance.

2106 (8) Multi-model ensemble with diverse preprocessing: cluster trained models by their error
2107 patterns and combine their outputs. Apply different preprocessing pipelines (raw,
2108 background-removed, edge maps, thresholded) to the input, denoise with clustered sub-
2109 ensembles, and fuse predictions for robustness across noise characteristics.

2110 (9) Enhanced augmentation and target refinement: simulate realistic scanning artifacts by
2111 applying 3D text deformations, lighting variations, and shadows to clean images. Identify
2112 and compensate for systematic artifacts in the provided 'clean' targets (e.g., single-
2113 pixel halos) by either modeling them as noise or adjusting ground-truth labels during
2114 training.

2113 **And solution drafts generated in this iteration are:**

2114 Pipeline 1: ResNet-34 Encoder U-Net with Multi-Scale Edge & Total-Variation Loss

2115 Overview:

2116 A robust baseline using a pretrained ResNet-34 backbone as a U-Net encoder fused with a light-
2117 weight decoder. Combines L1 loss, Sobel edge loss at multiple scales, and a total-
2118 variation regularizer to preserve text strokes while smoothing background noise. Mixed
2119 precision training and sliding-window inference ensure the entire pipeline runs in ~45
2120 min on an A6000.

2121 1. Data Preparation

2122 - Read all train noisy/clean PNGs, normalize to [0,1].

2123 - Extract on-the-fly 256*256 patches: random crop + random horizontal/vertical flips + 90
2124 degree rotations.

2125 - 90/10 split by file indices (seed=42). Use batch size 8-16.

2126 2. Model Architecture

2127 - Encoder: torchvision.models.resnet34(pretrained=True), first conv modified to 1->64
2128 channels.

2129 - Decoder: four upsampling stages (ConvTranspose2d + Conv2d+BN+ReLU) mirroring ResNet blocks,
2130 with skip-connections from encoder layers.

2131 - Final conv 64->1 + Sigmoid.

2132 3. Loss Function

2133 Let \hat{y} and y be predictions and targets.

2134 - L1Loss(\hat{y} , y)

2135 - Edge loss: L1 between Sobel(\hat{y}) and Sobel(y) at both full resolution and half resolution
2136 (downsample by 2).

2137 - TV: $\lambda \cdot TV(\hat{y})$ where $TV = \text{mean}(|\nabla_x \hat{y}| + |\nabla_y \hat{y}|)$.

2138 Total loss = $\alpha \cdot L1 + \beta \cdot \text{Edge_full} + \gamma \cdot \text{Edge_half} + \delta \cdot TV$, e.g. $\alpha=1.0$, β
2139 $=0.5$, $\gamma=0.25$, $\delta=1e-5$.

2140 4. Optimization

2141 - Optimizer: AdamW(lr=1e-3, weight_decay=1e-4).

2142 - Scheduler: CosineAnnealingLR(T_max=25).

2143 - Mixed precision via torch.cuda.amp.

2144 - Early stopping on validation RMSE (patience=5).

2145 5. Inference & Submission

2146 - Perform sliding-window inference on each test image with 256*256 patches, stride=128.

2147 - Average overlapping patches.

2148 - Clamp outputs to [0,1], write submission.csv with id=image_row_col.

2149 Compute budget: ~20 min train + ~5 min inference.

2150 Pipeline 2: Laplacian-Pyramid Multi-Scale Residual U-Net with Pyramid Loss

2151 Overview:

2152 A novel pyramid-domain network that decomposes images into multi-scale Laplacian bands,
2153 denoises each band via shared-weight residual blocks, and merges them back. Multi-level
2154 L1 losses focus the model on both coarse structures and fine text details. Efficient and
2155 fully end-to-end in PyTorch, training finishes in ~50 min on an A6000.

2156 1. Data & Augmentation

2157 - Same data split and patch sampling (256*256, flips, rotations).

2158 2. Laplacian Pyramid Transform (LPT)

2159 - On the fly, for each patch: create 2-level Gaussian pyramid using avg-pooling (scale
2160 1->0.5->0.25), then compute Laplacian bands $L0=(I-\text{upsample}(I/2))$, $L1=(I/2-\text{upsample}(I/4))$,
2161 and a residual low band $L2=I/4$.

2162 - Stack [L0,L1,L2] as 3 input channels.

2163 3. Network

2164 - Encoder-decoder U-Net with 3-channel input and 3-channel output, 4 down/upsampling stages,
2165 64->128->256->512->256->128->64 channel widths.

2166 - Each block: two 3*3 Conv+BN+ReLU and a residual skip connection.

2167 4. Reconstruction

2168 - Output is 3 residual corrections R0,R1,R2 to each band.

2169 - Reconstruct clean patch:


```

2160     a0 = L0 + R0; a1 = L1 + R1; a2 = L2 + R2;
2161     I = a0 + up(a1) + up(up(a2)).
2162
2163 5. Loss Function
2164 - Per-band L1 losses:  $w_0 \cdot L_1(a_0, H_0) + w_1 \cdot L_1(a_1, H_1) + w_2 \cdot L_1(a_2, H_2)$ , where  $H_k$  are ground-truth
2165   bands and  $w_0=1.0$ ,  $w_1=0.5$ ,  $w_2=0.25$ .
2166 - Pixel L1:  $\lambda_{\text{pix}} \cdot L_1(I, I_{\text{gt}})$  with  $\lambda_{\text{pix}}=0.1$ .
2167
2168 6. Training
2169 - AdamW(lr=5e-4, weight_decay=1e-5), one-cycle scheduler (max_lr=5e-4, epochs=30).
2170 - Batch size 8, mixed precision.
2171 - Early stopping on val RMSE (patience=6).
2172
2173 7. Inference
2174 - On full test images, run sliding-window with 256*256 patches, stride=128.
2175 - For each patch, apply pyramid -> model -> reconstruct -> accumulate and average.
2176 - Clamp to [0,1], write submission.csv.
2177
2178 Compute budget: ~25 min training + ~8 min inference.

```

The best code proposed in the second iteration is:

```

2175 1 import os, random, csv
2176 2 import numpy as np
2177 3 from glob import glob
2178 4 from PIL import Image
2179 5
2180 6 import torch
2181 7 import torch.nn as nn
2182 8 import torch.nn.functional as F
2183 9 from torch.utils.data import Dataset, DataLoader
2184 10 from torchvision import transforms, models as tv_models
2185 11 from torch.cuda.amp import autocast, GradScaler
2186 12
2187 13 # 1) Reproducibility
2188 14 SEED = 42
2189 15 random.seed(SEED)
2190 16 np.random.seed(SEED)
2191 17 torch.manual_seed(SEED)
2192 18 torch.cuda.manual_seed_all(SEED)
2193 19
2194 20 DEVICE = torch.device("cuda" if torch.cuda.is_available() else "cpu")
2195 21
2196 22 # 2) Paths
2197 23 TRAIN_NOISY = "./input/train"
2198 24 TRAIN_CLEAN = "./input/train_cleaned"
2199 25 TEST_DIR = "./input/test"
2200 26 WORK_DIR = "./working0"
2201 27 os.makedirs(WORK_DIR, exist_ok=True)
2202 28
2203 29 # 3) Hyperparameters
2204 30 PATCH_SIZE = 256
2205 31 STRIDE = 128
2206 32 BATCH_SIZE = 12
2207 33 LR, WD = 1e-3, 1e-4
2208 34 MAX_EPOCHS = 50
2209 35 PATIENCE = 10
2210 36 T_MAX = 50 # for LR scheduler
2211 37
2212 38 # 4) Loss-weight constants (from tuning)
2213 39 w1, w2, w3, w4 = 1.0, 0.5, 0.25, 1e-5
2214 40 lambda_aux = 0.4394633936788146
2215 41 lambda_mse = 0.1312037280884873
2216 42 lambda_ssim = 0.031198904067240532
2217 43 lambda_ssim2 = lambda_ssim / 2
2218 44
2219 45 # 5) Dataset + augmentations
2220 46 class OCRDataset(Dataset):
2221 47

```

```

2214 48 def __init__(self, noisy_list, clean_list, ps, train):
2215 49     self.noisy, self.clean = noisy_list, clean_list
2216 50     self.ps, self.train = ps, train
2217 51     self.to_tensor = transforms.ToTensor()
2218 52     self.aug = transforms.Compose(
2219 53         [
2220 54             transforms.RandomChoice(
2221 55                 [
2222 56                     transforms.RandomHorizontalFlip(1.0),
2223 57                     transforms.RandomVerticalFlip(1.0),
2224 58                     transforms.RandomRotation(90),
2225 59                     transforms.RandomRotation(180),
2226 60                     transforms.RandomRotation(270),
2227 61                 ]
2228 62             ),
2229 63             transforms.RandomApply([transforms.GaussianBlur(3, (0.1, 2.0)
2230 64                                     ), p=0.3]),
2231 65             transforms.RandomApply([transforms.RandomAdjustSharpness(2.0)
2232 66                                     ], p=0.3),
2233 67         ]
2234 68     )
2235 69
2236 70 def __len__(self):
2237 71     return len(self.noisy)
2238 72
2239 73 def __getitem__(self, i):
2240 74     n = Image.open(self.noisy[i]).convert("L")
2241 75     c = Image.open(self.clean[i]).convert("L")
2242 76     w, h = n.size
2243 77     # pad
2244 78     if w < self.ps or h < self.ps:
2245 79         pad = (0, 0, max(0, self.ps - w), max(0, self.ps - h))
2246 80         n = transforms.functional.pad(n, pad, fill=255)
2247 81         c = transforms.functional.pad(c, pad, fill=255)
2248 82         w, h = n.size
2249 83     # crop
2250 84     if self.train:
2251 85         x = random.randint(0, w - self.ps)
2252 86         y = random.randint(0, h - self.ps)
2253 87     else:
2254 88         x = (w - self.ps) // 2
2255 89         y = (h - self.ps) // 2
2256 90     n = n.crop((x, y, x + self.ps, y + self.ps))
2257 91     c = c.crop((x, y, x + self.ps, y + self.ps))
2258 92     if self.train and random.random() < 0.5:
2259 93         n = self.aug(n)
2260 94         c = self.aug(c)
2261 95     return self.to_tensor(n), self.to_tensor(c)
2262 96
2263 97 # 6) Prepare train/val split
2264 98 noisy_files = sorted(glob(f"{TRAIN_NOISY}/*.png"))
2265 99 clean_files = [f"{TRAIN_CLEAN}/{os.path.basename(x)}" for x in
2266 100               noisy_files]
2267 101 N = len(noisy_files)
2268 102 idx = list(range(N))
2269 103 random.shuffle(idx)
2270 104 ntr = int(0.9 * N)
2271 105 tr_idx, va_idx = idx[:ntr], idx[ntr:]
2272 106 train_noisy = [noisy_files[i] for i in tr_idx]
2273 107 train_clean = [clean_files[i] for i in tr_idx]
2274 108 val_noisy = [noisy_files[i] for i in va_idx]
2275 109 val_clean = [clean_files[i] for i in va_idx]
2276 110 train_ds = OCRDataset(train_noisy, train_clean, PATCH_SIZE, train=True)

```

```

2268 val_ds = OCRDataset(val_noisy, val_clean, PATCH_SIZE, train=False)
2269 train_loader = DataLoader(
2270     train_ds, batch_size=BATCH_SIZE, shuffle=True, num_workers=4,
2271     pin_memory=True
2272 )
2273 val_loader = DataLoader(
2274     val_ds, batch_size=BATCH_SIZE, shuffle=False, num_workers=4,
2275     pin_memory=True
2276 )
2277 # 7) Sobel, TV, SSIM helpers
2278 sob_x = (
2279     torch.tensor([[1, 0, -1], [2, 0, -2], [1, 0, -1]], dtype=torch.float32
2280 )
2281     .view(1, 1, 3, 3)
2282     .to(DEVICE)
2283 )
2284 sob_y = sob_x.transpose(2, 3)
2285
2286 def sobel(x):
2287     gx = F.conv2d(x, sob_x, padding=1)
2288     gy = F.conv2d(x, sob_y, padding=1)
2289     return torch.sqrt(gx * gx + gy * gy + 1e-6)
2290
2291 def total_variation(x):
2292     dh = (x[:, :, 1:, :] - x[:, :, :-1, :]).abs().mean()
2293     dw = (x[:, :, :, 1:] - x[:, :, :, :-1]).abs().mean()
2294     return dh + dw
2295
2296 def ssim_map(a, b, C1=0.01**2, C2=0.03**2):
2297     mu_a = F.avg_pool2d(a, 3, 1, 1)
2298     mu_b = F.avg_pool2d(b, 3, 1, 1)
2299     sa = F.avg_pool2d(a * a, 3, 1, 1) - mu_a * mu_a
2300     sb = F.avg_pool2d(b * b, 3, 1, 1) - mu_b * mu_b
2301     sab = F.avg_pool2d(a * b, 3, 1, 1) - mu_a * mu_b
2302     num = (2 * mu_a * mu_b + C1) * (2 * sab + C2)
2303     den = (mu_a * mu_a + mu_b * mu_b + C1) * (sa + sb + C2)
2304     return num / (den + 1e-8)
2305
2306 def ssim_loss(a, b):
2307     return 1.0 - ssim_map(a, b).mean()
2308
2309 # 8) loss_terms
2310 l1_loss = nn.L1Loss()
2311 mse_loss = nn.MSELoss()
2312
2313 def loss_terms(pred, target):
2314     L1v = l1_loss(pred, target)
2315     MSEv = mse_loss(pred, target)
2316     Ef = l1_loss(sobel(pred), sobel(target))
2317     p2, t2 = F.avg_pool2d(pred, 2), F.avg_pool2d(target, 2)
2318     Eh = l1_loss(sobel(p2), sobel(t2))
2319     TVv = total_variation(pred)
2320     return L1v, MSEv, Ef, Eh, TVv
2321
2322 # 9) Model w/ deep supervision
2323 class ResUNetDS(nn.Module):
2324     def __init__(self):

```

```

2322 172 super().__init__()
2323 173 r34 = tv_models.resnet34(pretrained=True)
2324 174 self.enc0 = nn.Conv2d(1, 64, 7, 2, 3, bias=False)
2325 175 self.enc0.weight.data = r34.conv1.weight.data.mean(dim=1, keepdim=
2326 176 True)
2327 177 self.bn0, self.relu0, self.pool0 = r34.bn1, r34.relu, r34.maxpool
2328 178 self.enc1, self.enc2 = r34.layer1, r34.layer2
2329 179 self.enc3, self.enc4 = r34.layer3, r34.layer4

2330 180 def up(i, o):
2331 181     return nn.ConvTranspose2d(i, o, 2, 2)

2332 182 def cb(i, o):
2333 183     return nn.Sequential(
2334 184         nn.Conv2d(i, o, 3, 1, 1, bias=False),
2335 185         nn.BatchNorm2d(o),
2336 186         nn.ReLU(inplace=True),
2337 187         nn.Conv2d(o, o, 3, 1, 1, bias=False),
2338 188         nn.BatchNorm2d(o),
2339 189         nn.ReLU(inplace=True),
2340 190     )

2341 191 self.up4, self.dec4 = up(512, 256), cb(256 + 256, 256)
2342 192 self.up3, self.dec3 = up(256, 128), cb(128 + 128, 128)
2343 193 self.up2, self.dec2 = up(128, 64), cb(64 + 64, 64)
2344 194 self.aux_up, self.aux_out = up(64, 64), nn.Conv2d(64, 1, 1)
2345 195 self.up1, self.dec1 = up(64, 64), cb(64 + 64, 64)
2346 196 self.up0, self.outc = up(64, 64), nn.Conv2d(64, 1, 1)
2347 197 self.sig = nn.Sigmoid()

2348 200 def forward(self, x):
2349 201     x0 = self.relu0(self.bn0(self.enc0(x)))
2350 202     x1 = self.pool0(x0)
2351 203     x2 = self.enc1(x1)
2352 204     x3 = self.enc2(x2)
2353 205     x4 = self.enc3(x3)
2354 206     x5 = self.enc4(x4)

2355 207     d4 = self.dec4(torch.cat([self.up4(x5), x4], dim=1))
2356 208     d3 = self.dec3(torch.cat([self.up3(d4), x3], dim=1))
2357 209     d2 = self.dec2(torch.cat([self.up2(d3), x2], dim=1))
2358 210     aux = self.sig(self.aux_out(self.aux_up(d2)))
2359 211     d1 = self.dec1(torch.cat([self.up1(d2), x0], dim=1))
2360 212     main = self.sig(self.outc(self.up0(d1)))
2361 213     return main, aux

2362 214
2363 215 model = ResUNetDS().to(DEVICE)
2364 216
2365 217 # 10) Optimizer, scheduler, scaler
2366 218 optimizer = torch.optim.AdamW(model.parameters(), lr=LR, weight_decay=WD)
2367 219 scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, T_max=
2368 220 T_MAX)
2369 221 scaler = GradScaler()

2370 222
2371 223 # 11) Training + snapshot saving
2372 224 best_rmse = float("inf")
2373 225 patience = 0
2374 226 snap_epochs = set([10, 20, 30, 40, 50])

2375 227
2376 228 for epoch in range(1, MAX_EPOCHS + 1):
2377 229     model.train()
2378 230     train_loss = 0.0
2379 231     for noisy_img, clean_img in train_loader:
2380 232         noisy_img, clean_img = noisy_img.to(DEVICE), clean_img.to(DEVICE)

```

```

2376 optimizer.zero_grad()
2377 with autocast():
2378     main_pred, aux_pred = model(noisy_img)
2379     Llv, MSEv, Ef, Eh, TVv = loss_terms(main_pred, clean_img)
2380     s1 = ssim_loss(main_pred, clean_img)
2381     p2, t2 = F.avg_pool2d(main_pred, 2), F.avg_pool2d(clean_img, 2)
2382     s2 = ssim_loss(p2, t2)
2383     main_loss = (
2384         w1 * Llv
2385         + lambda_mse * MSEv
2386         + w2 * Ef
2387         + w3 * Eh
2388         + w4 * TVv
2389         + lambda_ssim * s1
2390         + lambda_ssim2 * s2
2391     )
2392     aux_up = F.interpolate(
2393         aux_pred,
2394         size=clean_img.shape[-2:],
2395         mode="bilinear",
2396         align_corners=False,
2397     )
2398     La, Ma, Ea, Eh2, TVa = loss_terms(aux_up, clean_img)
2399     sa = ssim_loss(aux_up, clean_img)
2400     pa, ca = F.avg_pool2d(aux_up, 2), F.avg_pool2d(clean_img, 2)
2401     sa2 = ssim_loss(pa, ca)
2402     aux_loss = (
2403         w1 * La
2404         + lambda_mse * Ma
2405         + w2 * Ea
2406         + w3 * Eh2
2407         + w4 * TVa
2408         + lambda_ssim * sa
2409         + lambda_ssim2 * sa2
2410     )
2411     loss = main_loss + lambda_aux * aux_loss
2412     scaler.scale(loss).backward()
2413     scaler.step(optimizer)
2414     scaler.update()
2415     train_loss += loss.item()
2416 scheduler.step()
2417
2418 # validation
2419 model.eval()
2420 se, count = 0.0, 0
2421 with torch.no_grad():
2422     for noisy_img, clean_img in val_loader:
2423         noisy_img, clean_img = noisy_img.to(DEVICE), clean_img.to(DEVICE)
2424         with autocast():
2425             pred, _ = model(noisy_img)
2426             se += ((pred - clean_img) ** 2).sum().item()
2427             count += pred.numel()
2428 val_rmse = np.sqrt(se / count)
2429 print(
2430     f"Epoch {epoch}: TrainLoss={train_loss/len(train_loader):.4f},
2431     ValRMSE={val_rmse:.6f}"
2432 )
2433
2434 # best + snapshot
2435 if val_rmse < best_rmse:
2436     best_rmse = val_rmse
2437     torch.save(model.state_dict(), os.path.join(WORK_DIR, "best_full.
2438         pth"))
2439     patience = 0

```

```

2430
2431297
2432298
2433299
2434300
2435301
2436302
2437303
2438304
2439305
2440306
2441307
2442308
2443309
2444310
2445311
2446312
2447313
2448314
2449315
2450316
2451317
2452318
2453319
2454320
2455321
2456322
2457323
2458324
2459325
2460326
2461327
2462328
2463329
2464330
2465331
2466332
2467333
2468334
2469335
2470336
2471337
2472338
2473339
2474340
2475341
2476342
2477343
2478344
2479345
2480346
2481347
2482348
2483349
2484350
2485351
2486352
2487353
2488354
2489355
2490356
2491357
2492358

else:
    patience += 1
if epoch in snap_epochs:
    torch.save(model.state_dict(), os.path.join(WORK_DIR, f"snap_{epoch}
        ").pth))
if patience >= PATIENCE:
    print("Early stopping.")
    break

print("Best validation RMSE:", best_rmse)

# 12) Ensemble load
ckpts = ["best_full.pth"] + sorted(
    [f for f in os.listdir(WORK_DIR) if f.startswith("snap_")],
    key=lambda x: int(x.split("_")[1].split(".")[0]),
)[-2:]
ensemble_nets = []
for ck in ckpts:
    net = ResUNetDS().to(DEVICE)
    net.load_state_dict(torch.load(os.path.join(WORK_DIR, ck)))
    net.eval()
    ensemble_nets.append(net)

# 13) Sliding-window ensemble inference
def ensemble_infer(img_arr):
    h, w = img_arr.shape
    inp = torch.from_numpy(img_arr / 255.0).unsqueeze(0).unsqueeze(0).to(
        DEVICE)
    ph = (PATCH_SIZE - h % STRIDE) % STRIDE
    pw = (PATCH_SIZE - w % STRIDE) % STRIDE
    inp = F.pad(inp, (0, pw, 0, ph), mode="reflect")
    _, _, H, W = inp.shape
    out = torch.zeros_like(inp)
    wt = torch.zeros_like(inp)
    for y in range(0, H - PATCH_SIZE + 1, STRIDE):
        for x in range(0, W - PATCH_SIZE + 1, STRIDE):
            patch = inp[:, :, y : y + PATCH_SIZE, x : x + PATCH_SIZE]
            preds = []
            with torch.no_grad(), autocast():
                for net in ensemble_nets:
                    p, _ = net(patch)
                    preds.append(p)
            avg_p = torch.stack(preds, 0).mean(0)
            out[:, :, y : y + PATCH_SIZE, x : x + PATCH_SIZE] += avg_p
            wt[:, :, y : y + PATCH_SIZE, x : x + PATCH_SIZE] += 1.0
    out = out / wt
    out = out[:, :, :h, :w]
    return out.detach().cpu().numpy().squeeze()

# 14) Write submission.csv
submission_path = os.path.join(WORK_DIR, "submission.csv")
with open(submission_path, "w", newline="") as f:
    writer = csv.writer(f)
    writer.writerow(["id", "value"])
    for tf in sorted(
        glob(f"{TEST_DIR}/*.png"), key=lambda x: int(os.path.basename(x).
            split(".")[0])
    ):
        img_id = os.path.basename(tf).split(".")[0]
        img = np.array(Image.open(tf).convert("L"), dtype=np.float32)
        den = ensemble_infer(img)
        H, W = den.shape
        for i in range(H):

```

```
2484359         for j in range(W):
2485360             writer.writerow([f"{img_id}_{i+1}_{j+1}", f"{den[i,j]:.6f}"])
2486361 print("Submission saved to", submission_path)
```

```
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
```