CoMind: Towards Community-Driven Agents for Machine Learning Engineering

Anonymous Author(s)

Affiliation Address email

Abstract

Large language model-based machine learning (ML) agents have shown great promise in automating ML research. 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 novel agent that excels at exchanging insights and developing novel solutions within a community context. CoMind achieves state-of-the-art performance on MLE-Live and outperforms 79.2% human competitors on average across four ongoing Kaggle competitions.

12 1 Introduction

3

6

8

10

11

20

21

23 24

25

26

27

28

29

30

31

32

33

35

Large language model (LLM)-based agents have shown remarkable potential in automating complex reasoning and decision-making tasks, ranging from software engineering (Jimenez et al., 2023a; Xia et al., 2025) and mathematical problem solving (OpenAI, 2024; Ren et al., 2025; Li et al., 2025) to scientific research exploration (Romera-Paredes et al., 2024; Yamada et al., 2025; Sun et al., 2025; Feng et al., 2025). Among these domains, machine learning engineering (MLE) remains a particularly impactful yet challenging application area, requiring design, implementation, and evaluation of high-performing models across diverse data science tasks.

Recent advances have introduced LLM agents capable of autonomously developing machine learning pipelines for Kaggle-style competitions (Chan et al., 2025). For example, MLAB (Huang et al., 2024) adopts a ReAct-style (Yao et al., 2023) agent for structured decision-making across tasks. AIDE (Jiang et al., 2025) leverages tree-based exploration for improved efficiency, and AutoKaggle (Li et al., 2024) introduces a multi-agent system with skill specialization. These agents have made progress toward end-to-end automation of MLE.

However, existing systems typically operate in isolation, relying solely on internal memory and trial-and-error exploration while ignoring a critical component of real-world scientific practice — community knowledge sharing. In real data science competitions and research workflows, participants frequently learn from public discussions, shared notebooks, and community insights. Such collective knowledge significantly enhances solution quality and innovation. Current agents, due to the inability to engage with this dynamic external context, often converge to repetitive strategies and plateau in performance. Therefore, we are motivated to explore the following critical research question:

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

To address this question, we introduce **MLE-Live**, a novel live evaluation framework simulating a Kaggle-style research community. Unlike prior benchmarks, MLE-Live includes time-stamped

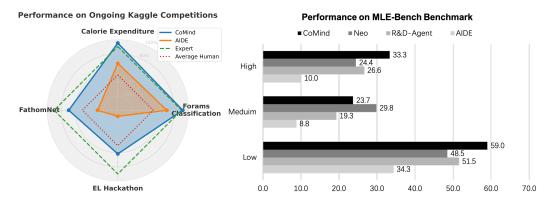


Figure 1: **Left:** CoMind's win rates on four 4 ongoing Kaggle competitions. **Right:** CoMind achieves state-of-the-art performance on the MLE-Bench competitions, measured by *Any Medal* score.

public discussions and shared code artifacts available before competition deadlines – resources that human participants routinely leverage. This setup allows us to rigorously evaluate an agent's ability to use community knowledge in a realistic, temporally grounded setting. In addition, MLE-Live supports both offline evaluation on past competitions and online evaluation on ongoing competitions, enabling comprehensive assessment across static and dynamic scenarios.

Building upon this framework, we propose **CoMind**, a new LLM-based MLE agent that systematically explores diverse ideas, iteratively refines solutions, and selectively incorporates external knowledge. CoMind maintains an evolving idea pool and constructs multiple distinct solution drafts in parallel. It dynamically focuses on one draft at a time, enabling efficient implementation without prompt overflow while preserving technical accuracy. Inspired by human brainstorming, this design balances exploratory breadth with practical depth.

We evaluate CoMind in both previous and ongoing data science competitions. For evaluation on previous competitions, CoMind is tested on the MLE-Live benchmark, which includes 75 past Kaggle competitions on MLE-Bench. CoMind achieves state-of-the-art performance on the leaderboard, significantly outperforming prior agents such as AIDE and R&D-Agent. For evaluation on ongoing competitions, we deploy CoMind on four ongoing Kaggle competitions, where it outperforms 79.2% human competitors on average (Figure 1), demonstrating strong real-world practicality. Human evaluations further confirm that CoMind generates more sophisticated and longer code, reflecting deeper reasoning and better integration of novel insights.

55 In summary, our contributions are:

- MLE-Live: A live evaluation framework simulating community-driven machine learning research, including shared discussions and code for realistic agent benchmarking.
- **CoMind**: A novel LLM-based agent that excels at leveraging collective knowledge and iterative idea exploration, achieving medal-level performance in real Kaggle competitions.
- Agent design innovations: We propose a new iterative and parallel exploration mechanism that
 enables continuous knowledge accumulation, improves diversity, and overcomes LLM context
 window limitations.

2 Related Work

56

57

58

59

60

61

62

63

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., 2023b).

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 86 capabilities. MLPerf (Mattson et al., 2020) assesses system-level performance, including training 87 speed and energy efficiency. To evaluate end-to-end ML workflows, MLAB (Huang et al., 2024) 88 tests the capabilities of LLM-based agents across 13 ML tasks. MLE-Bench (Chan et al., 2025) and 89 DSBench (Jing et al., 2025) further extends to about 75 Kaggle competitions covering tasks such 90 as preprocessing, modeling, and evaluation. However, these benchmarks typically evaluate agents 91 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 93 agents' ability to engage with and benefit from shared knowledge — while ensuring that resource 94 access remains fair and realistic.

3 MLE-Live

115

116 117

Existing benchmarks typically evaluate ML agents in static, isolated settings, where agents work independently without interacting with other participants or leveraging community insights. This contrasts sharply with real-world machine learning workflows, particularly on platforms like Kaggle, where collaboration, public sharing, and discussion are essential drivers of innovation.

To enable more realistic and comprehensive evaluation of agents in community-driven research settings, we introduce MLE-Live, a live evaluation framework built upon Kaggle-style competitions.

MLE-Live extends the MLE-Bench benchmark (Chan et al., 2025) by incorporating simulated community interactions that mirror how human participants access and utilize shared knowledge during competitions.

Each competition in MLE-Live includes: (i) Task description: Background information, task specifications, evaluation metrics, and dataset structure scraped directly from the original Kaggle competition pages. (ii) Competition dataset: A cleaned train-test split, constructed when necessary to account for unavailable or partial test data after competition closure. (iii) Submission grader: An evaluation script that mimics Kaggle's scoring mechanism. (iv) Leaderboard: A ranking system to reflect solution quality and progress.

Simulated Community Environment. Beyond the standard competition setup, MLE-Live introduces a community simulation system that mimics the collaborative ecosystem of real Kaggle competitions. Specifically, we collect and curate:

- Shared discussions, including strategy brainstorming, model diagnostics.
- Public kernels, including end-to-end solutions and code snippets which were posted before the
 official competition deadline.

These artifacts reflect the auxiliary resources human participants would naturally reference, making MLE-Live a richer and more authentic testbed for ML agents.

In total, we collected 12,951 discussions and 15,733 kernels across 75 Kaggle competitions from MLE-Bench.

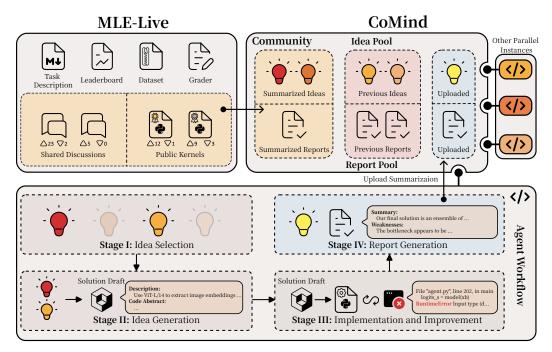


Figure 2: Overview of MLE-Live and the workflow of CoMind. CoMind simulates a community and operates iteratively through four stages: *Idea Selection, Idea Generation, Implementation and Improvement*, and *Report Generation*. Multiple agents on the same task share the community knowledge base, new reports will be added and visible to others in subsequent iterations.

Metadata and Quality Signals. Each resource in MLE-Live is augmented with critical metadata to help agents and evaluators prioritize relevant, high-quality content:

- Vote count: Community preference indicator; highly voted items often contain well-structured insights.
- **Public score:** Automatically computed by Kaggle on the public test split, indicating kernel performance.
- Author tier: A qualitative marker of the contributor's expertise, ranging from Novice to Grandmaster.

Importantly, all included content was published before competition deadlines, ensuring a faithful simulation of real-time knowledge sharing without post-hoc leakage. This design makes MLE-Live a controlled yet realistic benchmark to assess how well agents can leverage collective intelligence during the research process.

4 CoMind

122

123

124

125

126

127

128

129

130

131

132

133

134

We introduce **CoMind**, a community-augmented large language model (LLM) agent designed to 135 automate machine learning (ML) engineering in an iterative, collaborative setting. Figure 2 is a 136 overview of CoMind workflows. Inspired by the workflow of human practitioners on platforms 137 like Kaggle, CoMind operates in a loop that mirrors how experts read community posts, form new 138 ideas, experiment, and share results. The system operates in iterative cycles, each consisting of four 139 stages: Idea Selection, Idea Generation, Implementation and Improvement, and Report Generation. 140 To support cumulative progress, CoMind maintains two central repositories: an *idea pool* containing 141 abstracted insights derived from community content and prior iterations, and a report pool containing 142 finalized solution reports with associated code, evaluations, and analyses. Additionally, we extend 143 the setting to support multi-agent collaboration through shared insight exchange. These components 144 facilitate both intra-agent memory and inter-agent communication in the multi-agent deployments.

146 4.1 Agent Workflow

Stage I: Idea Selection. At the beginning of each iteration, CoMind accesses the idea pool, which contains curated concepts and strategies distilled from previous solutions, public kernels and forum discussions. By utilizing the report pool as a guidance of performance and relevance assessment of entries in the idea pool, CoMind ranks and filters these entries to identify a subset of ideas most promising for the current task. This process mimics how human participants explore collective wisdom before forming new hypotheses and experimenting.

Stage II: Idea Generation. Based on the selected ideas and additional context from the report pool, which contains detailed descriptions of previous solution implementations and their empirical performance, CoMind generates a high-level *solution draft*. This draft synthesizes new strategies by recombining or extending the selected ideas. Importantly, it is designed to avoid simple replication, thereby ensuring conceptual diversity and promoting exploratory breadth. This stage reflects the human capacity for abstraction and innovation, where participants generalize from past work to create novel solution blueprints.

Stage III: Implementation and Improvement. Based on the generated solution draft, CoMind 160 initiates a ReAct-style loop to implement, validate and refine the pipeline. In this stage, the CoMind 161 iteratively issues code snippets, executes them within a constrained and isolated runtime environment, 162 observes feedback (e.g., validation metrics, error logs), and updates its implementation accordingly. 163 This loop continues for a fixed time budget, allowing CoMind to incrementally debug and optimize its 164 solution through trial and error. Notably, the agent's contextual input during this stage is deliberately 165 restricted to include only the problem statement, the specific solution draft, and execution feedback, excluding direct access to the broader ideal pool and report pool. This ensures that the CoMind 167 develops the solution path independently, preserving experimental modularity while preventing the 168 underlying explosion of context length. 169

Stage IV: Report Generation. Upon convergence or budget exhaustion, CoMind compiles a comprehensive report for the solution draft, which consist of: (1) a clear description of the proposed method; (2) an analysis of each major component; (3) quantitative performance results; and (4) an assessment of limitations and future directions. The resulting report is then posted back to the simulated community by added to the report pool, making it available to other agents in future iterations. This mirrors how real users document and share their final solutions.

4.2 Parallel Agents with Shared Insight

Beyond a single-agent loop, CoMind also supports a collaborative multi-agent setting. Multiple agents operate in parallel on the same task, each with access to the shared community knowledge base. As agents generate new reports, these are added to the pool and can be read by others in subsequent iterations. This allows agents to build upon each other's ideas, fostering community-driven exploration and collective improvement.

182 5 Benchmark Evaluation

5.1 Setup

176

183

Task Selection. Based on MLE-Live evaluation framework, we evaluate our agent on 75 Kaggle 184 competitions on MLE-Bench. Using the MLE-Live framework, CoMind has access to shared 185 186 discussions and public kernels published on the competition websites before the competition deadline. To validate CoMind under realistic conditions, we further evaluate CoMind on four ongoing 187 Kaggle competitions: el-hackathon-2025, fathomnet-2025, playground-series-s5e5 and 188 forams-classification-2025. These competitions span diverse domains, including tabular learn-189 ing, image classification and 3D object classification. Rather than approximating the official scoring 190 locally, we directly submit CoMind's generated submission.csv files to the Kaggle platform, so that all reported ranks reflect genuine, live leaderboard positions. Notably, fathomnet-2025 is part 192 of the Fine-Grained Visual Categorization (FGVC12) workshop at the CVPR Conference. Unlike typical Kaggle competitions, a panel of experts will review the top entries based not only on scores

Table 1: **Any Medal** (%) **scores on MLE-Bench competitions.** Best results in each column are highlighted in bold. Baseline results are from the official 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

but also on the methodological descriptions. Although the competition offers no monetary prize, it serves as a high-profile venue for academic and practical contributions to marine biodiversity research.

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. The *Implementation and Improvement* stage of CoMind is limited to a maximum of 20 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 leaderboard 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 MLAgentBench, a general tool-calling/plan-act baseline; and Neo (https://heyneo.so/), a close-sourced multi-agent system for autonomous ML engineering.

5.2 Results

Table 1 compares CoMind with baseline methods on 75 MLE-Bench competitions. CoMind achieves state-of-the-art performance with an *Any Medal* rate of 36.00%, significantly outper-forming open-source competitors such as R&D-Agent (submitted on 2025-08-15) and surpassing the closed-source multi-agent system Neo. Appendix C provides a detailed case study on denoising-dirty-documents.

https://github.com/openai/mle-bench

Table 2: Authentic scores and top-percentile ranks of CoMind and AIDE on ongoing Kaggle competitions. "Higher better" marks whether a larger score is better for that competition. "Top %" is the percentile rank on competition leaderboard (lower is better). In playground-series-s5e5 and fathomnet-2025, lower scores are better.

		CoMind		AIDE	
Competition	Score higher better	Score	Top %	Score	Top %
playground-series-s5e5	×	0.5673	5.1%	0.5772	33.8%
forams-classification-2025	✓	0.7645	8.3%	0.6041	30.6%
el-hackathon-2025	✓	0.5837	38.4%	0.1140	91.5%
fathomnet-2025	×	2.81	30.6%	3.71	71.4%

On the four evaluated ongoing competitions CoMind's standings are: playground-series-s5e5 (#120 out of 2,338); forams-classification-2025 (#4 out of 48); el-hackathon-2025 (#128 out of 333); fathomnet-2025 (#15 out of 47). Details including authentic scores and win rates per task are provided in Table 2. These authentic results demonstrate CoMind's capability to tackle a variety of problem domains and achieve competitive performance in live, evolving ML workflows. In particular, our success in the CVPR-affiliated fathomnet-2025 challenge highlights CoMind's potential to contribute meaningfully not only to industrial applications but also to scientific and interdisciplinary research communities.

236 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—selected by highest community votes. This version, 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 R. In this variant, CoMind operates without access to any external community resources. It starts with empty idea and report pools 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.
- **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.

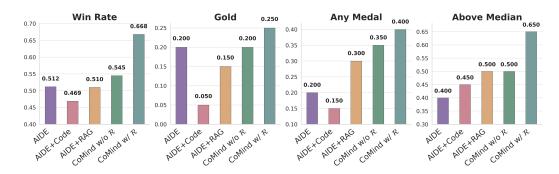


Figure 3: **Performance of CoMind and other baselines on 20 competitions from MLE-Bench-Lite.** *Valid Submission* is the ratio of submissions meeting format requirements and validation criteria. *Win Rate* is the percentage of human competitors outperformed by the agent. *Any Medal*, is the proportion 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 3: 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. Notably, CoMind demonstrated superior performance in Image Classification, Text Classification, Audio Classification and Image To Image.

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

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.

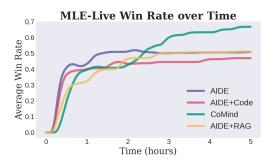
6.2 Results

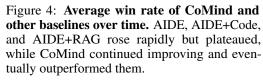
Figure 3 shows the results. Our key findings are as follows: (i) CoMind consistently outperforms all baselines across every metric. (ii) Among the AIDE variants, AIDE+RAG outperforms AIDE+Code, and both surpass the original AIDE on most metrics, demonstrating the benefits of integrating community knowledge. 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.

Task Categories Table 3 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





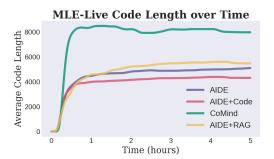


Figure 5: Average code length (character count) of valid solutions over time. CoMind maintained a substantially longer code length, suggesting more complex logic and richer optimization techniques.

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 4 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 5 illustrates the average code length during the entire competition. CoMind consistently generates significantly longer and more complex code, while other baselines begin with simpler implementations and introduce only incremental modifications. Appendix A offers a comparative analysis across code complexity metrics and task categories. Notably, CoMind's solutions for Image Regression and Audio Classification are nearly twice as long as those of other baselines. Additionally, solutions from CoMind are, on average, 55.4% longer than those produced by AIDE.

8 Conclusion

We introduce **MLE-Live**, a new framework for evaluating machine learning agents in realistic, community-driven environments. By simulating the collaborative dynamics of Kaggle competitions with shared discussions and public code, MLE-Live enables a more faithful assessment of agents' ability to leverage collective knowledge. Building upon this framework, we propose **CoMind**, an LLM-based agent that iteratively selects, synthesizes, and implements ideas using both internal reasoning and external insights. CoMind consistently outperforms prior methods on MLE-Live benchmark and four ongoing Kaggle competitions, demonstrating the value of community awareness and iterative exploration in research automation. In addition, MLE-Live lays the groundwork for future studies in collaborative AI systems where agents not only learn from data, but also from each other.

Limitations and Future Work. Currently, CoMind supports only report-level interactions. Expanding the agent's action space to include commenting, question-asking, or sharing datasets and models is a promising next step. In addition, 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.

18 References

- AI-Researcher (2025). Ai-researcher: Fully-automated scientific discovery with llm agents. Accessed: 2025-05-15.
- Boiko, D. A., MacKnight, R., Kline, B., and Gomes, G. (2023). Autonomous chemical research with large language models. *Nature*, 624:570 578.
- Chan, J. S., Chowdhury, N., Jaffe, O., Aung, J., Sherburn, D., Mays, E., Starace, G., Liu, K., Maksin, L., Patwardhan, T., Madry, A., and Weng, L. (2025). MLE-bench: Evaluating machine learning
- agents on machine learning engineering. In The Thirteenth International Conference on Learning
 Representations.
- Chi, Y., Lin, Y., Hong, S., Pan, D., Fei, Y., Mei, G., Liu, B., Pang, T., Kwok, J., Zhang, C., Liu, B., and Wu, C. (2024). Sela: Tree-search enhanced llm agents for automated machine learning. *ArXiv*,
- abs/2410.17238.
- Erickson, N., Mueller, J., Shirkov, A., Zhang, H., Larroy, P., Li, M., and Smola, A. (2020). Autogluontabular: Robust and accurate automl for structured data. *arXiv* preprint arXiv:2003.06505.
- Feng, S., Sun, W., Li, S., Talwalkar, A., and Yang, Y. (2025). A comprehensive evaluation of contemporary ml-based solvers for combinatorial optimization. *ArXiv*, abs/2505.16952.
- Feurer, M., Eggensperger, K., Falkner, S., Lindauer, M., and Hutter, F. (2022). Auto-sklearn 2.0:
 Hands-free automl via meta-learning. *Journal of Machine Learning Research*, 23(261):1–61.
- Grosnit, A., Maraval, A. M., Doran, J., Paolo, G., Thomas, A., Beevi, R. S. H. N., Gonzalez, J.,
- Khandelwal, K., Iacobacci, I., Benechehab, A., Cherkaoui, H., Hili, Y. A. E., Shao, K., Hao, J.,
- Yao, J., Kégl, B., Bou-Ammar, H., and Wang, J. (2024). Large language models orchestrating
- structured reasoning achieve kaggle grandmaster level. ArXiv, abs/2411.03562.
- Guo, S., Deng, C., Wen, Y., Chen, H., Chang, Y., and Wang, J. (2024). Ds-agent: Automated data science by empowering large language models with case-based reasoning. *ArXiv*, abs/2402.17453.
- Hollmann, N., Müller, S. G., and Hutter, F. (2023). Large language models for automated data science: Introducing caafe for context-aware automated feature engineering. In *Neural Information*
- 344 Processing Systems.
- Hong, S., Lin, Y., Liu, B., Wu, B., Li, D., Chen, J., Zhang, J., Wang, J., Zhang, L., Zhuge, M., Guo,
- 346 T., Zhou, T., Tao, W., Wang, W., Tang, X., Lu, X., Liang, X., Fei, Y., Cheng, Y., Gou, Z., Xu,
- Z., Wu, C., Zhang, L., Yang, M., and Zheng, X. (2024). Data interpreter: An Ilm agent for data
- science. *ArXiv*, abs/2402.18679.
- Hong, S., Zheng, X., Chen, J. P., Cheng, Y., Zhang, C., Wang, Z., Yau, S. K. S., Lin, Z. H., Zhou, L.,
- Ran, C., Xiao, L., and Wu, C. (2023). Metagpt: Meta programming for multi-agent collaborative
- s51 framework. *ArXiv*, abs/2308.00352.
- Huang, Q., Vora, J., Liang, P., and Leskovec, J. (2024). MLAgentbench: Evaluating language
- agents on machine learning experimentation. In Forty-first International Conference on Machine
- 354 Learning.
- Jiang, Z., Schmidt, D., Srikanth, D., Xu, D., Kaplan, I., Jacenko, D., and Wu, Y. (2025). Aide:
 Ai-driven exploration in the space of code.
- Jimenez, C. E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., and Narasimhan, K. (2023a). Swebench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*.
- Jimenez, C. E., Yang, J., Wettig, A., Yao, S., Pei, K., Press, O., and Narasimhan, K. (2023b). Swe-bench: Can language models resolve real-world github issues? *ArXiv*, abs/2310.06770.
- Jing, L., Huang, Z., Wang, X., Yao, W., Yu, W., Ma, K., Zhang, H., Du, X., and Yu, D. (2025).
- DSBench: How far are data science agents from becoming data science experts? In *The Thirteenth*
- 363 International Conference on Learning Representations.

- Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., and Talwalkar, A. (2018). Hyperband: A novel
 bandit-based approach to hyperparameter optimization. *Journal of Machine Learning Research*,
 18(185):1–52.
- Li, S., Marwah, T., Shen, J., Sun, W., Risteski, A., Yang, Y., and Talwalkar, A. (2025). Codepde: An inference framework for llm-driven pde solver generation. *arXiv preprint arXiv:2505.08783*.
- Li, Z., Zang, Q., Ma, D., Guo, J., Zheng, T., Liu, M., Niu, X., Wang, Y., Yang, J., Liu, J., et al. (2024).

 Autokaggle: A multi-agent framework for autonomous data science competitions. *arXiv preprint arXiv:2410.20424*.
- Liu, H., Simonyan, K., and Yang, Y. (2019). DARTS: Differentiable architecture search. In *International Conference on Learning Representations*.
- Liu, Z., Cai, Y., Zhu, X., Zheng, Y., Chen, R., Wen, Y., Wang, Y., Weinan, E., and Chen, S. (2025).

 Ml-master: Towards ai-for-ai via integration of exploration and reasoning. *ArXiv*, abs/2506.16499.
- Mattson, P., Cheng, C., Diamos, G., Coleman, C., Micikevicius, P., Patterson, D., Tang, H., Wei, G.-Y., Bailis, P., Bittorf, V., et al. (2020). Mlperf training benchmark. *Proceedings of Machine Learning and Systems*, 2:336–349.
- OpenAI (2024). Learning to reason with llms.
- OpenAI (2025). Introducing openai o3 and o4-mini.
- Ren, Z. Z., Shao, Z., Song, J., Xin, H., Wang, H., Zhao, W., Zhang, L., Fu, Z., Zhu, Q., Yang, D., Wu, Z. F., Gou, Z., Ma, S., Tang, H., Liu, Y., Gao, W., Guo, D., and Ruan, C. (2025). Deepseek-prover-v2: Advancing formal mathematical reasoning via reinforcement learning for subgoal decomposition.
- Romera-Paredes, B., Barekatain, M., Novikov, A., Balog, M., Kumar, M. P., Dupont, E., Ruiz, F. J., Ellenberg, J. S., Wang, P., Fawzi, O., et al. (2024). Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475.
- Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Zettlemoyer, L., Cancedda, N., and Scialom, T. (2023). Toolformer: Language models can teach themselves to use tools.
- Shen, Y., Song, K., Tan, X., Li, D., Lu, W., and Zhuang, Y. (2023). Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face.
- Sun, W., Feng, S., Li, S., and Yang, Y. (2025). Co-bench: Benchmarking language model agents in algorithm search for combinatorial optimization. *ArXiv*, abs/2504.04310.
- Thornton, C., Hutter, F., Hoos, H. H., and Leyton-Brown, K. (2013). 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*, pages 847–855.
- Trirat, P., Jeong, W., and Hwang, S. J. (2024). Automl-agent: A multi-agent llm framework for full-pipeline automl. *ArXiv*, abs/2410.02958.
- Wang, C., Wu, Q., Weimer, M., and Zhu, E. (2021). Flaml: A fast and lightweight automl library.
 Proceedings of Machine Learning and Systems, 3:434–447.
- Wang, L., Yang, N., Huang, X., Yang, L., Majumder, R., and Wei, F. (2024). Multilingual e5 text embeddings: A technical report. *arXiv preprint arXiv:2402.05672*.
- Wang, X., Li, B., Song, Y., Xu, F. F., Tang, X., Zhuge, M., Pan, J., Song, Y., Li, B., Singh, J., Tran, H. H., Li, F., Ma, R., Zheng, M., Qian, B., Shao, Y., Muennighoff, N., Zhang, Y., Hui, B., Lin,
- H. H., Li, F., Ma, R., Zheng, M., Qian, B., Shao, Y., Muennighoff, N., Zhang, Y., Hui, B., Lin, J., Brennan, R., Peng, H., Ji, H., and Neubig, G. (2025). Openhands: An open platform for AI
- software developers as generalist agents. In *The Thirteenth International Conference on Learning*
- 407 Representations.
- Xia, C. S., Deng, Y., Dunn, S., and Zhang, L. (2025). Demystifying llm-based software engineering
 agents. Proceedings of the ACM on Software Engineering, 2(FSE):801–824.

- Xie, T., Zhang, D., Chen, J., Li, X., Zhao, S., Cao, R., Hua, T. J., Cheng, Z., Shin, D., Lei, F., Liu, Y.,
 Xu, Y., Zhou, S., Savarese, S., Xiong, C., Zhong, V., and Yu, T. (2024). 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.
- Yamada, Y., Lange, R. T., Lu, C., Hu, S., Lu, C., Foerster, J., Clune, J., and Ha, D. (2025). The ai scientist-v2: Workshop-level automated scientific discovery via agentic tree search. *arXiv preprint arXiv:2504.08066*.
- Yang, X., Yang, X., Fang, S., Xian, B., Li, Y., Wang, J., Xu, M., Pan, H., Hong, X., Liu, W., Shen, Y., Chen, W., and Bian, J. (2025). R&d-agent: Automating data-driven ai solution building through llm-powered automated research, development, and evolution. *ArXiv*, abs/2505.14738.
- Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K. R., and Cao, Y. (2023). React:
 Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*.
- Zhou, S., Xu, F. F., Zhu, H., Zhou, X., Lo, R., Sridhar, A., Cheng, X., Ou, T., Bisk, Y., Fried, D.,
 Alon, U., and Neubig, G. (2024). Webarena: A realistic web environment for building autonomous
 agents. In *The Twelfth International Conference on Learning Representations*.

A Additional Analysis on Code Complexity

427

430

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

Table 4: Code Complexity and Quality Metrics by Task Category.

Category	Metric	CoMind	AIDE	AIDE+RAG	AIDE+Code
	CC	1.68	1.59	1.93	1.29
Image Classification	Pylint Score	7.43	9.06	8.90	8.92
	Volume	330.88	143.26	84.20	175.88
	Difficulty	4.95	2.90	2.32	2.59
	Effort	1960.22	507.06	286.31	725.59
	SLOC	198.25	133.50	120.88	115.71
	Comment Lines	15.62	12.88	13.75	14.43
	Code Length	7638.40	4624.30	4701.30	5192.10
	CC	3.58	4.28	2.00	0.00
	Pylint Score	8.82	9.09	8.89	9.26
	Volume	286.38	384.07	47.68	29.25
T . C1	Difficulty	3.76	3.94	1.25	1.31
Text Classification	Effort	1183.11	2332.22	61.56	35.16
	SLOC	181.67	133.00	141.00	69.50
	Comment Lines	14.67	15.33	14.00	13.50
	Code Length	6974.70	3094.50	5920.50	5629.30
	CC	2.00	0.00	0.00	0.00
	Pylint Score	7.92	9.11	9.49	8.86
	Volume	718.63	244.20	115.95	227.48
	Difficulty	7.39	6.46	3.19	6.38
Audio Classification	Effort	5308.07	1577.11	369.58	1451.30
	SLOC	256.00	82.00	92.00	72.00
	Comment Lines	20.00	11.00	16.00	16.00
	Code Length	9449.00	3508.00	4151.00	3352.00
	CC	4.38	2.25	22.33	15.75
	Pylint Score	8.58	9.04	9.14	8.51
	Volume	492.55	52.33	390.46	324.00
	Difficulty	3.87	2.14	5.26	3.68
Seq2Seq	Effort	1935.02	140.58	2083.84	1686.74
	SLOC	184.50	63.50	222.50	147.50
	Comment Lines	22.50	13.00	23.00	19.50
	Code Length	6925.50	5649.50	8357.50	2728.50
	CC	2.78	1.62	2.38	0.25
	Pylint Score	8.65	8.96	8.87	9.31
	Volume	1264.61	856.12	815.29	435.46
	Difficulty	7.37	4.83	6.05	3.69
Tabular	Effort	10 808.93	6163.62	5564.22	2001.06
	SLOC	218.75	139.75	147.50	93.50
	Comment Lines	18.25	133.75 14.75	15.25	10.50
	Code Length	8570.00	3534.00	6064.00	5759.80
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 Code Length	$26.00 \\ 8800.00$	$8.00 \\ 5231.00$	23.00 4815.00	$13.00 \\ 6671.00$
	-				
	CC	1.68	2.00	2.40	2.00

Category	Metric	CoMind	AIDE	AIDE+RAG	AIDE+Code
	Pylint Score	8.62	8.75	8.80	8.89
	Volume	1310.92	241.08	70.32	72.00
	Difficulty	8.75	3.88	2.18	2.73
	Effort	11466.58	934.17	153.43	196.36
	SLOC	267.00	145.00	116.00	133.00
	Comment Lines	36.00	15.00	12.00	12.00
	Code Length	10991.00	4841.00	4655.00	5614.00

431

B Prompts and Responses for CoMind

- This section provides some examples of prompts and responses in CoMind, including **Idea Selection**,
- Idea Generation, Implementation and Improvement and Report Generation.

435 B.1 Idea Selection

- The idea pool is initialized with curated strategies distilled from public kernels and forum discussions
- before the first iteration. Since the key innovation of CoMind lies in the simulation of the community,
- we adopt a relatively simple implementation for the distillation, where CoMind only collects and
- analyzes top-k voted or ranked (with best public score) kernels and discussions.

Prompt for Strategy Distillation of Public Kernels

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>

440

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.

442

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'].

443 444

445

Once the idea pool is initialized, CoMind enters the main iteration. CoMinds then ranks and filters all entries in the idea pool, following the prompt below.

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>

446

Response Template of Idea Filtering and Reconstruction

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'].

447

448

B.2 Idea Generation

Based on previous ideas and reports. CoMind will first enrich the idea pool by designing and generating other promising strategies for the competition.

Prompt for CoMind 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".
- 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 I>
- <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.

After brainstorming completes, CoMind synthesizes existing strategies and ideas into several high-level *solution drafts*.

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.

452

453

- 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

456

Response Template for Solution Draft Synthesis

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

457

458 B.3 Implementation

In this stage, all previously generated solution drafts will be distributed to multiple parallel sub-agents.

In each instance, CoMind chats with LLM in multiple turns and initiates a ReAct-style loop.

Prompts for Sub-Agent 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>

Reminders

- 1. Read the pipeline and task description carefully.
- 2. YOUR CODE MUST PRODUCE SUBMISSON AT ./working-agent_id/submission.csv, THIS IS EXTREMELY IMPORTANT
- 3. There is one A6000 gpu available for you, maximize your use of computing resources. You can use large batchsizes.
- 4. All the provided input data are stored in ./input directory.
- 5. You can use the ./working-agent_id directory to store any temporary files that your code needs to create.
- 6. Include at least one comment explaining your code. NO PARTS OF THE CODE SHOULD BE SKIPPED OR OMITTED, don't terminate before finishing the script. Even if your proposed code is a minor change, don't omit any sections that overlap with the previous code.
- 7. Remember, your ultimate goal is to Obtain best score on this competition.
- 8. Your code should **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. This is important because we will pick your best code based on this metric.
- 10. We suggest you to test your code at a small scale and print necessary information before utilizing full dataset to get familiar with the data structure and avoid potential format errors.
- 11. Time limit per run is 1 hour. Your code will be killed if timeout.

12. Begin by summarizing your understanding of the task, and then propose your first code.

Response Format You should follow the following format:

Response Template for Sub-Agent Implementation

objective of this implementation and suggestions for output evaluation key technical considerations

expected running time (you should ensure that the code will finish within 1 hour)

```
''python
your code here
```

462

A Python environment will be setup and execute the code automatically. After the execution, a summary LLM collects all standard output and determines whether the execution runs successfully.

Prompt for Execution Result Summarization

Introduction You are a Kaggle grandmaster attending a competition. You have written code to solve this task and now need to evaluate the output of the code execution. You should determine if there were any bugs as well as report the empirical findings. Include essential information about the result, including warnings, errors, and the final metric.

Code <*Python code generated by the agent>*

Goals and Explanation <explanation of the code>

Execution Output <terminal output>

465

Response Template for Execution Result Summarization

is_bug true if the output log shows that the execution failed or has some bug, otherwise false

summary write a short summary (4-5 sentences) describing the empirical findings **output_abs** select representative segments of the output log and mark the remainder as ellipses

metric If the code ran successfully and produced submission.csv on full test set (i.e. not dummy or partial), report the value of the final validation metric. Otherwise, leave it null. is_lower_better true if the metric should be minimized (i.e. a lower metric value is better, such as with MSE), false if the metric should be maximized (i.e. a higher metric value is better, such as with accuracy)

466 467

468

CoMind progressively optimizes the solution through trail and error. After the execution result is collected and summarized, it will try to revise the code by notifying the LLM:

Prompt for Consequent Code Revisions

Remaining steps: <remaining_steps>; Remaining time: <remaining_time> seconds
I ran your code and summarized the execution result: <summary>

Now, please choose your next action and propose code using the same response format as before. Remember, output a self-contained code, no part of it should be omitted. Keep the final validation metric same as the metric mentioned in the task description.

- A) Fix runtime errors (if any)
- B) Do hyperparameter tuning
- C) Include ideas that were not implemented yet
- D) Add possible improvements

E) Run on a larger scale (moderately increase training epochs, etc.). You should refer to the previous execution time we reported. Remember, your code will be killed if timeout.

470

71 **B.4 Report Generation**

Each sub-agent of CoMind compiles a comprehensive report and submits it to the report pool when time expires.

Prompt for Report Compilation

Please summarize the results and submit a comprehensive report.

474

Response Template for 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)

<vour 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)

<vour rationale>

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

475

480

476 C Case Study: Denoising Dirty Documents

477 C.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.
 - C.1.1 Example of Public Kernel

```
481
482
    A simple feed-forward neural network that denoises one pixel at a time
483
484
    import numpy as np
485
    import theano
486
    import theano.tensor as T
48%
    import cv2
488
    import os
489
    import itertools
490
    theano.config.floatX = 'float32'
```

```
493
    def load_image(path):
494
        return cv2.imread(path, cv2.IMREAD_GRAYSCALE)
495
496
    def feature_matrix(img):
497
        """Converts a grayscale image to a feature matrix
498
499
        The output value has shape (<number of pixels>, <number of features>)
500
501)
        # select all the pixels in a square around the target pixel as features
502
503
        window = (5, 5)
        nbrs = [cv2.getRectSubPix(img, window, (y, x)).ravel()
504
                for x, y in itertools.product(range(img.shape[0]), range(img.shape[1]))]
505
506
        # add some more possibly relevant numbers as features
507
        median5 = cv2.medianBlur(img, 5).ravel()
508
        median25 = cv2.medianBlur(img, 25).ravel()
509
        grad = np.abs(cv2.Sobel(img, cv2.CV_16S, 1, 1, ksize=3).ravel())
510
        div = np.abs(cv2.Sobel(img, cv2.CV_16S, 2, 2, ksize=3).ravel())
5110
513
    ... (omitted) ...
513
513
        # for fname in os.listdir('../input/test/'):
515
        for fname in ['1.png']:
516
            test_image = load_image(os.path.join('.../input/test', fname))
517
            test_x = feature_matrix(test_image)
518
519
            y_pred, = predict(test_x)
520
521
            output = y_pred.reshape(test_image.shape)*255.0
522
            cv2.imwrite('original_' + fname, test_image)
523
            cv2.imwrite('cleaned_' + fname, output)
524
525
526
    if __name__ == '__main__':
527
        main()
528
```

C.1.2 Example of Discussion

530

553

554

556

```
531
     # Edge Diffraction in train_cleaned data
532
     (Lance <TIER: N/A>) I'm studying the pixels in train_cleaned data.&nbsp; I attached a colorized blow-up
           version of part of the image train_cleaned/45.png.   The yellow pixels are any pixels that
533
           were not pure white ( != 0xFF gray scale) in image 45.png, the green was pure white (0xFF).
535
        So you see what looks like an edge diffraction line lining the outer edge of all the letters.
        Okay, maybe I got something wrong in my code.  Can anyone confirm this edge diffraction thing
536
              in the train_cleaned data, as for example the first word in train_cleaned/45.png (There). 
537
              You need to make the non-white (byte != OxFF) pixels all a more contrasting color or you may not
538
539
540
        You guessing that the clean png files were at some point scanned in using some kind of optical
541
              scanning machine which added these edge diffraction lines when the light diffracts off the edge
              of the black ink character.
543
            (omitted) ...
        + (Rangel Dokov <TIER: MASTER>) Yes, there is some noise, which doesn't look like it should be there
545
               in the clean set... I ran a test setting everything whiter that 0xF5 to 0xFF and the RMSE was
              0.005, which should be an upper bound on the effects from the halos. This will likely be large
             enough to make the top of the leaderboard a game of luck, but since this is just a playground
548
              competition I'm not terribly worried about it.
```

C.2 Idea Selection

In our experiment settings, CoMind only accesses top-10 voted discussions and kernels and ignores the rest. Upon completion of this process, 7 ideas and 10 pipelines are generated. Below is an excerpt of the ideas and public pipelines.

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

```
(2) Apply a Waifu2x-inspired deep convolutional neural network with gradually increasing filter counts (e. g., 1 -> 32 -> 64 -> 128 -> 256 -> 512 -> 1) and LeakyReLU activations for effective denoising
557
558
         Carefully initialize convolutional weights (e.g., stdv = sqrt(2/(kW*kH*nOutputPlane))) and use
559
           LeakyReLU to improve model convergence and performance
560
         Ensemble multiple models with different input preprocessing: combine outputs from a pure CNN,
561
           background-removed images, edge maps, and thresholded inputs to capture diverse noise characteristics
562
     (5) Augment training data to simulate real-world 3D deformations and shadows on text, not just 2D noise,
563
564
           to better match test-time artifacts
     (6) Account for systematic artifacts in 'clean' training data (e.g., single-pixel halos) by treating them
565
566
           as noise or adjusting targets accordingly during training
567
     Public pipeline (0): - Pipeline: A simple feed-forward neural network that denoises one pixel at a time (
568
           Theano).
569
       - Feature engineering: for each pixel extract a 5*5 window of gray values (neighbors), 5*5 median blur,
570
             25*25 median blur, Sobel gradient and second-order derivative magnitudes, stack into a feature
571
             vector. Normalize features to [0,1].
572
       - Model architecture: two-layer MLP; hidden layer size N_HIDDEN=10, tanh activation, output layer with
573
             custom activation clip(x+0.5,0,1).
574
       - Training: MSE cost, stochastic gradient descent with learning rate 0.1, batch size 20, epochs 100.
575
             Validation on one image (3.png) at each epoch by RMSE.
576
       - Prediction: apply same feature_matrix to test images, predict pixel values, reshape to full image,
577
             write out cleaned PNGs.
     - Code abstract:
578
         def feature_matrix(img):
579
             window=(5,5)
580
581
             nbrs=[cv2.getRectSubPix(img,window,(y,x)).ravel()
582
                   for x,y in itertools.product(range(img.shape[0]),range(img.shape[1]))]
583
             median5=cv2.medianBlur(img,5).ravel()
             median25=cv2.medianBlur(img,25).ravel()
585
             grad=np.abs(cv2.Sobel(img,cv2.CV_16S,1,1,ksize=3).ravel())
             div=np.abs(cv2.Sobel(img,cv2.CV_16S,2,2,ksize=3).ravel())
586
587
             misc=np.vstack((median5,median25,grad,div)).T
588
             features=np.hstack((nbrs,misc))
             return (features/255.).astype('float32')
589
590
591
         class Model(object):
592
             def __init__(...):
593
                \verb|self.layer1=Layer(...,n_in=...,n_out=N_HIDDEN,activation=T.tanh)|\\
                self.layer2=Layer(...,n_in=N_HIDDEN,n_out=n_out,
594
595
                                  activation=lambda x: T.clip(x+0.5,0,1))
596
            def cost(self,y): return T.mean((self.output-y)**2)
597
598
      ----- PIPELINE SEPARATOR -----
599
     Public pipeline (1): - Pipeline: Matching image backgrounds in R (no ML model).
       - Reads test PNGs in batches of 12 images.
600
         Flattens each into vectors of size 258*540, stacks as columns.
601
         For each pixel location, takes the maximum value across images as an estimate of background.
602
         Writes out background images as PNG.
603
604
      - Code abstract:
         for(i in 1:4) {
605
           matches = seq(1,205,by=12) + (i-1)*3
606
           rawData=matrix(0,258*540,length(matches))
607
608
           for(j in seq_along(matches)){
             imgY=readPNG(file.path(testDir,paste0(matches[j],'.png')))
609
             rawData[,j]=as.vector(imgY[1:258,1:540])
610
611
           background=matrix(apply(rawData,1,max),258,540)
612
           writePNG(background, paste0('background',matches[j],'.png'))
613
614
615
616
      ----- PIPELINE SEPARATOR -----
     Public pipeline (2): - Pipeline: Pixel-wise Random Forest regression (Python, chunk size=1e6).
617
618
       - Feature engineering: pad image by mean value (padding=1); extract 3*3 neighborhood per pixel, flatten
619
             as features.
620
       - Training data: load all train noisy images, compute features via joblib parallel (n_jobs=128), load
621
             targets as flattened clean pixel intensities/255.
622
       - Model: sklearn.ensemble.RandomForestRegressor(warm_start=True, n_jobs=-1). Incrementally add one
623
             estimator at a time: split training rows into CHUNKSIZE=1e6 slices, in each slice increase
624
             n\_estimators by 1 and fit on that slice.
625
       - Prediction: extract test features similarly, generate index strings "image_row_col", predict pixel
626
            values, write submission CSV.
     - Code abstract:
627
628
         def get_padded(img, padding=1):
629
            padval=int(round(img.mean()))
630
              ... return padded
         def get_features_for_image(img,padding=1):
631
632
             padded=get_padded(img,padding)
633
             return np.vstack([padded[i:i+3,j:j+3].reshape(1,-1)
634
                              for i in range(rows) for j in range(cols)])
         def get_model(X,y):
636
```

```
637
             model=RandomForestRegressor(n_estimators=0,warm_start=True,n_jobs=-1)
638
             for start.end in slices:
639
                 model.set_params(n_estimators=model.n_estimators+1)
                 model.fit(X[start:end],y[start:end])
640
641
            return model
642
     ----- PIPELINE SEPARATOR -----
643
644
     ... (omitted) ..
      ----- PIPELINE SEPARATOR -----
645
     Public pipeline (9): - Pipeline: XGBoost with basic features in R.
646
       - Features: raw pixel value, kmeansThreshold binary, background_Removal feature.
- Sample 25k pixels, build DMatrix, 5-fold CV, early stop, select rounds, train xgb.
647
648
649
         Predict on test, clip, build id strings, write submission.csv.
650
     - Code abstract:
         dat=cbind(y,raw=img2vec(imgX),thresholded=kmeansThreshold(imgX),
651
                  backgroundRemoval=img2vec(background_Removal(imgX)))
652
653
         rows=sample(nrow(dat),25000)
         dtrain=xgb.DMatrix(dat[rows,-1],label=dat[rows,1])
654
         xgb.tab=xgb.cv(data=dtrain,nfold=5,early.stop.round=50)
655
656
         xgb.mod=xgboost(data=dtrain,nrounds=min.err.idx)
657
         yHat=predict(xgb.mod,newdata=as.matrix(x))
658
     After generating ideas and public pipelines, CoMind initiates the first iteration. It begins brainstorming
659
     to propose more original strategies.
660
661
     ===SOLUTION_PATH_1===
     Multi-Scale Wavelet-U-Net with Attention and Multi-Term Loss
662
     Description: decompose each image via discrete wavelet transform (DWT) into low-/high-frequency bands,
663
           process each band in a dual-branch U-Net enhanced with channel- and spatial-attention, then invert
665
           the wavelet to reconstruct. Train end-to-end with a weighted sum of pixel-wise RMSE, wavelet-band
           consistency, and edge preservation losses.
     - apply 2D DWT to noisy input -> obtain LL, LH, HL, HH sub-bands
667
     - feed LL into a "coarse" encoder branch, feed concatenated LH/HL/HH into a "detail" encoder branch
669
      - use a U-Net decoder to upsample each branch back to patch size, fuse via learned 1*1 convolutions
      - insert Convolutional Block Attention Modules (CBAM) after each encoder and decoder block
670
      - define loss = alpha*pixelRMSE(clean,output) + beta*bandRMSE(wavelet(clean),wavelet(output)) + gamma*
           edgeLoss(Sobel(clean),Sobel(output))
672
      - train on full images with AdamW and a cosine-annealing LR schedule
673
674
     ===SOLUTION_PATH_2===
675
676
     Stroke-Aware Conditional GAN with OCR-Guided Perceptual Loss
     Description: build a conditional GAN (generator = deep residual encoder-decoder, discriminator = PatchGAN)
677
            that not only minimizes pixel loss but also preserves text strokes-enforce a stroke-level loss via a
678
            pre-trained small CNN classifier that predicts presence/width of strokes. Add an OCR-based
679
           perceptual loss: feed predictions through a frozen OCR engine embedding and minimize distance to
680
           clean embedding.
681
     - implement generator as ResNet blocks + skip connections (64->128->256->128->64)
682
      - implement discriminator as 70*70 PatchGAN to focus on local texture
683
     - include L1 pixel loss + adversarial loss + stroke consistency loss (L1 between stroke-CNN features on
684
           clean vs. restored)
685
     - freeze a small text-structure CNN (trained on binary masks) to extract stroke features
686
     - run Tesseract (or lightweight OCR CNN) on restored vs. clean, extract penultimate-layer activations, add
687
            perceptual loss term
688
689
     - train with R1 gradient penalty and spectral normalization on discriminator
690
691
     ===SOLUTION PATH 3===
692
     Joint Dictionary Learning + Non-Local Patch Aggregation
     Description: learn paired dictionaries (D_noisy, D_clean) for small patches (e.g. 8*8) via coupled K-SVD.
693
694
           At test time, extract overlapping patches, compute sparse codes alpha under D_noisy via OMP,
           reconstruct clean patches = D_clean*alpha. Then, perform non-local means on the reconstructed patches
695
696
            to exploit self-similarity and average aggregates.
     - sample a large bank of noisy/clean patch pairs, initialize D_noisy, D_clean with DCT basis - run coupled K-SVD to minimize ||D_noisy*alpha - y_noisy|| + ||D_clean*alpha - y_clean|| w.r.t. D_noisy,
697
698
699
           D_clean,alpha
700
     - at test time, for each image patch y_noisy, compute alpha via Orthogonal Matching Pursuit (sparsity <= k)
701
702
     - reconstruct y_clean_est = D_clean*alpha for each patch
703
     - perform block-matching to find K nearest patches per reference patch (Euclid dist), stack them
704
     - aggregate reconstructed patches with non-local weights (e.g. Gaussian on reconstruction residual)
705
706
     ===SOLUTION PATH 4===
707
     Self-Supervised Blind Denoising via Noise2Void + Test-Time Adaptation
708
     Description: exploit purely noisy data-train a small U-Net with masked pixel prediction (Noise2Void) on
           each test image at inference (test-time training). The network learns to predict a pixel from its
709
710
           context, gradually adapting to local noise statistics, then you run a forward pass to obtain the
           cleaned image. No clean target needed.
711
712
     - define blind-spot or random masking scheme: mask 1% pixels per batch, replace with neighbors
     - build a lightweight CNN (e.g. 5 down/up blocks with skip connections) that predicts a full image
     - fine-tune this CNN on each test image for N_iter (e.g. 500 steps) using only masked L2 loss
```

```
    - use data augmentation (rotations, flips) on the single test image to diversify contexts
    - after adaptation, perform a clean forward pass without masking to get the denoised output
    - optionally ensemble outputs from multiple random initializations to reduce variance
```

These strategies are then added to the idea pool. To remove similar ideas and decompose overlapped ideas, a reconstruction is performed subsequently. 9 ideas are preserved after the reconstruction.

- (0) Use diverse ensembles by clustering models based on their error patterns and combining outputs from differently preprocessed inputs (e.g., raw, background-removed, edge maps, thresholded) to capture varied noise characteristics.
- 723 (1) Develop a sliding-window, patch-based model that takes an input region and predicts multiple denoised 724 and super-resolved pixels simultaneously for efficient document enhancement.
- 725 (2) Implement a Waifu2x-inspired deep convolutional network with progressively increasing filter counts
 726 and LeakyReLU activations, initialized using He normalization for robust convergence in denoising
 727 tasks.
- 728 (3) Augment the training set with simulated 3D deformations, lighting variations, and shadows on text to 729 better match real-world scanning artifacts.
- 730 (4) Identify and compensate for systematic artifacts in the "clean" training targets (e.g., single-pixel halos) by modeling them as noise or adjusting ground-truth labels during training.
- 732 (5) Build a wavelet-domain dual-branch U-Net: decompose inputs with 2D DWT into approximation and detail
 733 bands, encode separately, decode and fuse via learned 1*1 convolutions, integrate CBAM attention
 734 modules, and train end-to-end with a composite loss combining pixel RMSE, wavelet-band RMSE, and edge
 735 loss, optimized using AdamW with a cosine-annealing LR schedule.
- (6) Apply a GAN-based restoration framework: use a ResNet-based generator and 70*70 PatchGAN discriminator, train with a combination of L1 pixel loss, adversarial loss, stroke-consistency loss (via a frozen stroke-feature CNN), and perceptual OCR-feature loss, incorporating R1 gradient penalty and spectral normalization.
- 740 (7) Leverage coupled dictionary learning and non-local aggregation: learn noisy and clean dictionaries 741 simultaneously via coupled K-SVD, reconstruct patches using OMP, enhance self-similarity with block-742 matching, and fuse reconstructed patches with non-local weighting.
- 743 (8) Implement blind-spot test-time adaptation: randomly mask a small percentage of pixels per batch, train
 744 a lightweight CNN with skip connections on the single test image using masked L2 loss and
 745 augmentations, fine-tune per image for a few hundred steps, and optionally ensemble multiple
 746 initializations for variance reduction.

747 C.3 Idea Generation

b = self.bottleneck(e4)

Based on the idea pool, public pipelines and previous reports (currently unavailable), CoMind compiles 2 solution drafts.

```
Pipeline 1 (Baseline CNN - U-Net with L1 loss and simple augmentations)
750
751
752
     - A lightweight U-Net trained on 256*256 patches of noisy/clean document pairs.
753
     - Single-branch encoder-decoder with skip connections, ReLU activations, batch norm.
754
755
     - Loss: L1 pixel-wise loss (less sensitive to outliers than MSE), optional total variation regularization.
     - Optimizer: Adam; train on 1 A6000 within 1 hr.
756
757
758
     Data Preparation
     1. Read all noisy (input) and clean (target) train images, normalize intensities to [0,1].
759
     2. Extract random 256*256 patches (stride = 128) with matching noisy/clean pairs.
760
     3. Data augmentation: random horizontal/vertical flips, +-90 degree rotations.
761
     4. Create PyTorch DataLoader with batch_size=16 (fits 48 GB) for ~100 k patches.
762
763
764
     Network Architecture (PyTorch pseudocode)
765
766
     class UNet(nn.Module):
767
         def __init__(self):
768
            super().__init__()
769
             # Encoder
770
             self.enc1 = DoubleConv(1, 64)
             self.enc2 = Down(64,128)
771
772
             self.enc3 = Down(128.256)
             self.enc4 = Down(256,512)
773
774
             # Bottleneck
775
             self.bottleneck = Down(512,512)
776
             # Decoder
             self.up4 = Up(1024,256)
777
             self.up3 = Up(512,128)
778
779
             self.up2 = Up(256,64)
             self.up1 = Up(128,64)
780
781
             self.final = nn.Conv2d(64,1,kernel_size=1)
         def forward(self,x):
782
             e1=self.enc1(x)
783
             e2=self.enc2(e1)
784
785
             e3=self.enc3(e2)
             e4=self.enc4(e3)
```

```
788
             d4=self.up4(b,e4)
             d3=self.up3(d4,e3)
789
             d2=self.up2(d3,e2)
790
             d1=self.up1(d2,e1)
791
             return torch.sigmoid(self.final(d1))
792
     ...
793
     Helper modules:
794
     - DoubleConv = (Conv2d -> BatchNorm2d -> ReLU) *2
795
     - Down = (MaxPool2d -> DoubleConv)
796
     - Up = (ConvTranspose2d for upsampling -> concatenate skip -> DoubleConv)
797
798
799
     - Loss = L1Loss(output, target) + lambda*TV(output) (lambda=1e-5 for smoothness).
800
801
     - Optimizer = Adam(lr=1e-3, weight_decay=1e-5).
     - LR schedule: ReduceLROnPlateau(monitor=val_loss, factor=0.5, patience=5).
802
803
     - Train for up to 50 epochs; early-stop if val_loss stagnates.
     - Validation: hold out 10% patches to monitor RMSE.
804
805
806
     Inference
807
     - For each test image (e.g., 540*258), slide 256*256 window with stride=128, predict, and average
808
           overlapping outputs.
809
     - Threshold nothing; output raw [0,1] floats per pixel.
810
811
     Compute Budget
812
     - ~100 k patches, batch 16, ~6 k steps per epoch. On A6000: ~2-3 min/epoch => 50 epochs ~ 2 hr; with early
813
            stopping < 1 hr.
814
815
    Pipeline 2 (Advanced Wavelet U-Net with CBAM and Composite Loss)
816
817
818
     - Dual-branch U-Net operating in wavelet domain (Haar DWT) to explicitly denoise tonal and textural
819
     - CBAM (Convolutional Block Attention Modules) to adaptively weigh spatial/channel features.
820
821
     - Loss = alpha*L1_pixel + beta*L2_wavelet + gamma*EdgeLoss.
     - Optimizer = AdamW + CosineAnnealingLR.
822
823
824
     Data Preparation
     - Same as Pipeline 1 (256*256 patches + augmentations).
825
     - On-the-fly DWT: for each noisy patch, compute one-level Haar DWT -> yields approximation (A) and details
826
827
            (H,V,D).
828
829
     Network Architecture
830
     (implemented in PyTorch)
831
832
     class WaveletUNet(nn.Module):
833
         def __init__(self):
834
             super().__init__()
835
             # Shared CBAM-Res blocks for Approx and Detail branches
             self.encA1 = CBAMResBlock(1,64)
836
             self.encD1 = CBAMResBlock(3,64)
837
             self.pool = nn.MaxPool2d(2)
838
             self.encA2 = CBAMResBlock(64,128)
839
             self.encD2 = CBAMResBlock(64,128)
840
             # Bottleneck
841
             self.bottleneck = CBAMResBlock(256,256)
842
843
             # Decoder
             self.up2 = UpRes(256,128)
844
             self.up1 = UpRes(128,64)
self.final = nn.Conv2d(64,4,kernel_size=1)
845
846
         def forward(self, x):
847
             # x shape: [B,1,H,W]
848
             A, (H,V,D) = haar_dwt(x) # PyWavelets or custom conv
d0 = torch.cat([self.encA1(A), self.encD1(torch.cat([H,V,D],dim=1))],dim=1)
849
850
             p1 = self.pool(d0)
851
             d1 = torch.cat([self.encA2(p1[:, :128]), self.encD2(p1[:, 128:])],dim=1)
852
             b = self.bottleneck(self.pool(d1))
853
854
             u2 = self.up2(b,d1)
855
             u1 = self.up1(u2,d0)
             out = self.final(u1) # [B,4,H,W]
856
857
             # Reconstruct denoised x via inverse DWT
             A_hat, H_hat, V_hat, D_hat = out[:,0:1], out[:,1:2], out[:,2:3], out[:,3:4]
858
             return ihaar_dwt((A_hat, (H_hat,V_hat,D_hat)))
859
860
861
     Key building blocks:
     - CBAMResBlock: residual block with Conv-BN-ReLU->Conv-BN, plus CBAM (channel & spatial attention after
862
863
           second Conv).
864
     - UpRes: ConvTranspose2d upsample->concatenate skip->CBAMResBlock.
865
866
     - L1_pixel = L1(output, target).
867
    - L2_wavelet = L2(DWT(output) - DWT(target)).
```

```
- EdgeLoss = L1(Sobel(output), Sobel(target)).
- TotalLoss = alpha*L1_pixel + beta*L2_wavelet + gamma*EdgeLoss (alpha=1, beta=0.5, gamma=0.2).
869
870
871
872
     - Optimizer: AdamW(lr=1e-3, weight_decay=1e-4).
873
     - Scheduler: CosineAnnealingLR(T_max=50).
874
     - Batch size=8 (higher memory due to dual branches).
875
     - Epochs=50, early-stop by composite loss on val set.
876
877
     Inference
878
      - Same sliding-window approach as Pipeline 1.
879
880
881
     Compute Budget
      - CBAM and wavelet layers add ~1.5* compute but within A6000 48 GB: patch batch=8 yields ~250 steps/epoch,
882
883
            so 50 epochs ~2 hr. Use early-stop or T_max=30 to finish within 1 hr.
884
885
886
     References
887
     - CBAM: Woo et al., "CBAM: Convolutional Block Attention Module" (ECCV 2018).
888
     - Haar DWT: PyWavelets or custom conv filters.
889
     - Edge loss: Sobel via fixed conv kernels.
890
891
     Expected Benefits
892
     - Explicit noise suppression in high-frequency bands.
      - Attention modules focus on text strokes and edge restoration.
893
894
     - Composite loss balances global structure and fine details.
```

895 C.4 Implementation and Improvement

In our experiment, 2 parallel instances are invoked to implement the solution drafts respectively. Each instance has a quota of 20 steps, 3 hours in total and 1 hour per execution. After each execution, another summarization LLM is instructed to analyze the terminal output.

899 C.4.1 Example of Terminal Output

```
900
     Loaded 115 train images
901
     Sampled train/val patch coords: 80000 10000
     .../python3.11/site-packages/torch/optim/lr_scheduler.py:62: UserWarning: The verbose parameter is
902
903
           deprecated. Please use get_last_lr() to access the learning rate.
       warnings.warn(
904
905
     Traceback (most recent call last):
       File ".../agent_0.py", line 231, in <module>
906
907
         for xb, yb in train_loader:
908
       File ".../python3.11/site-packages/torch/utils/data/dataloader.py", line 708, in __next__
909
         data = self._next_data()
910
911
       File ".../python3.11/site-packages/torch/utils/data/dataloader.py", line 1480, in _next_data
         return self._process_data(data)
912
913
914
       File ".../python3.11/site-packages/torch/utils/data/dataloader.py", line 1505, in _process_data
915
916
       File ".../python3.11/site-packages/torch/_utils.py", line 733, in reraise
917
         raise exception
918
     ValueError: Caught ValueError in DataLoader worker process 0.
     Original Traceback (most recent call last):
       File ".../python3.11/site-packages/torch/utils/data/_utils/worker.py", line 349, in _worker_loop
920
         data = fetcher.fetch(index) # type: ignore[possibly-undefined]
922
923
       File ".../python3.11/site-packages/torch/utils/data/_utils/fetch.py", line 52, in fetch
         data = [self.dataset[idx] for idx in possibly_batched_index]
925
       File ".../python3.11/site-packages/torch/utils/data/_utils/fetch.py", line 52, in stcomp>
926
         data = [self.dataset[idx] for idx in possibly_batched_index]
927
928
929
       File ".../agent_0.py", line 121, in __getitem__
         n_patch = torch.from_numpy(n_patch).unsqueeze(0)
930
931
932
     ValueError: At least one stride in the given numpy array is negative, and tensors with negative strides
          are not currently supported. (You can probably work around this by making a copy of your array with
933
          array.copy().)
934
```

C.4.2 Example of Output Summary

```
936 Terminal output (truncated):
937 Loaded 115 train images
```

```
Sampled train/val patch coords: 80000 10000
938
     .../python3.11/site-packages/torch/optim/lr_scheduler.py:62: UserWarning: ...
939
     Traceback (most recent call last):
940
       File ".../agent_0.py", line 231, in <module>
941
942
         for xb, yb in train_loader:
943
     ValueError: At least one stride in the given numpy array is negative, and tensors with negative strides
944
           are not currently supported. (You can probably work around this by making a copy of your array with
945
946
           array.copy().)
     Execution summary:
947
     The training script fails during the first epoch's data loading. After loading 115 train images and
948
949
          sampling 80k train and 10k validation patches, the DataLoader worker crashes with a ValueError due to
950
           negative strides when converting flipped NumPy patches to tensors. As a result, the pipeline halts
951
           before any training or validation metrics are produced and no submission.csv is generated.
     Execution time: 3.81 seconds.
952
```

C.4.3 Example of Generated Snippet

954

The best code snippet generated in the first iteration of CoMind is:

```
import os, glob, random, math
956
     from PIL import Image
958
     import torch
     import torch.nn as nn
959
     import torch.nn.functional as F
960
     from torch.utils.data import Dataset, DataLoader
966
     from torchvision import transforms
962
963
964
     # -- 1) Haar DWT / inverse DWT --
965
966
     class HaarDWT(nn.Module):
         def __init__(self):
967
             super().__init__()
968
             wt = torch.zeros(4, 1, 2, 2)
969
             wt[0, 0, :, :] = 0.5
970
             wt[1, 0, 0, 0], wt[1, 0, 0, 1], wt[1, 0, 1, 0], wt[1, 0, 1, 1] = (
976
972
                 0.5,
                 -0.5.
973
974
                 0.5,
975
                 -0.5,
             )
976
             wt[2, 0, 0, 0], wt[2, 0, 0, 1], wt[2, 0, 1, 0], wt[2, 0, 1, 1] = (
977
978
                 0.5,
                 0.5,
979
980
                 -0.5,
                 -0.5,
981
982
983
             wt[3, 0, 0, 0], wt[3, 0, 0, 1], wt[3, 0, 1, 0], wt[3, 0, 1, 1] = (
                 0.5,
984
985
                 -0.5,
                 -0.5,
986
                 0.5,
987
             )
988
             self.register_buffer("weight", wt)
989
990
         def forward(self, x):
991
992
             return F.conv2d(x, self.weight, stride=2)
993
994
995
     class HaarIDWT(nn.Module):
996
         def __init__(self):
             super().__init__()
997
            wt = HaarDWT().weight.clone()
998
             self.conv = nn.ConvTranspose2d(4, 1, 2, stride=2, bias=False)
999
             self.conv.weight.data.copy_(wt)
1000
1001
             self.conv.weight.requires_grad_(False)
1002
```

```
def forward(self, coeffs):
1003
1004
             return self.conv(coeffs)
1005
1066
1067
     # -- 2) Sobel edge for EdgeLoss --
     class Sobel(nn.Module):
1068
1009
         def __init__(self):
             super().__init__()
1010
             kx = torch.tensor(
1056
             [[1, 0, -1], [2, 0, -2], [1, 0, -1]], dtype=torch.float32).view(1, 1, 3, 3)
1012
1013
             ky = torch.tensor(
1014
                 [[1, 2, 1], [0, 0, 0], [-1, -2, -1]], dtype=torch.float32
1065
             ).view(1, 1, 3, 3)
1016
             self.register_buffer("wx", kx)
1067
10168
             self.register_buffer("wy", ky)
1019
         def forward(self, x):
1020
             gx = F.conv2d(x, self.wx, padding=1)
1026
1022
             gy = F.conv2d(x, self.wy, padding=1)
             return torch.sqrt(gx * gx + gy * gy + 1e-6)
1023
1024
1025
     # -- 3) CBAM, ResBlock, UpRes, WaveletUNet --
1026
     class CBAM(nn.Module):
1027
         def __init__(self, c, r=16, k=7):
1028
             super().__init__()
1029
             self.mlp = nn.Sequential(
1030
                 nn.Linear(c, c // r, bias=False),
1036
                 nn.ReLU(inplace=True),
1032
                 nn.Linear(c // r, c, bias=False),
1033
             )
1034
             self.spatial = nn.Conv2d(2, 1, kernel_size=k, padding=k // 2, bias=False)
1035
1036
         def forward(self, x):
1037
             b, c, h, w = x.shape
1038
1039
             avg = F.adaptive_avg_pool2d(x, 1).view(b, c)
1040
             mx = F.adaptive_max_pool2d(x, 1).view(b, c)
             ca = torch.sigmoid(self.mlp(avg) + self.mlp(mx)).view(b, c, 1, 1)
1046
             x2 = x * ca
1042
             avgc = x2.mean(1, True)
1043
1044
             \maxc, _ = x2.\max(1, True)
             sa = torch.sigmoid(self.spatial(torch.cat([avgc, maxc], 1)))
1045
             return x2 * sa
1046
1047
1048
     class CBAMResBlock(nn.Module):
1049
         def __init__(self, inp, outp):
1056
1056
             super().__init__()
             self.conv1 = nn.Conv2d(inp, outp, 3, padding=1, bias=False)
1052
             self.bn1 = nn.BatchNorm2d(outp)
1053
             self.conv2 = nn.Conv2d(outp, outp, 3, padding=1, bias=False)
1054
             self.bn2 = nn.BatchNorm2d(outp)
1055
             self.cbam = CBAM(outp)
1056
1057
             self.skip = (
                 nn.Conv2d(inp, outp, 1, bias=False) if inp != outp else nn.Identity()
1058
1059
1066
             self.relu = nn.ReLU(inplace=True)
1066
1062
         def forward(self, x):
             idt = self.skip(x)
1063
             out = self.relu(self.bn1(self.conv1(x)))
1064
1065
             out = self.bn2(self.conv2(out))
             out = self.cbam(out) + idt
1066
             return self.relu(out)
1067
```

```
1068
1069
     class UpRes(nn.Module):
1070
         def __init__(self, ic, sc, oc):
1076
1072
             super().__init__()
             self.up = nn.ConvTranspose2d(ic, oc, 2, stride=2)
1073
1074
             self.block = CBAMResBlock(oc + sc, oc)
1075
         def forward(self, x, skip):
1076
1077
             x = self.up(x)
1078
             if x.shape[-2:] != skip.shape[-2:]:
                 x = F.interpolate(
1079
                     x, size=skip.shape[-2:], mode="bilinear", align_corners=False
1080
1081
             return self.block(torch.cat([x, skip], 1))
1082
1083
1084
     class WaveletUNet(nn.Module):
1085
         def __init__(self):
1086
1087
             super().__init__()
             self.dwt = HaarDWT()
1088
             self.idwt = HaarIDWT()
1089
             self.eA1 = CBAMResBlock(1, 64)
1090
             self.eD1 = CBAMResBlock(3, 64)
109₺
1092
             self.pool = nn.MaxPool2d(2)
             self.eA2 = CBAMResBlock(64, 128)
1093
             self.eD2 = CBAMResBlock(64, 128)
1094
             self.b = CBAMResBlock(256, 256)
1095
1096
             self.u2 = UpRes(256, 256, 128)
             self.u1 = UpRes(128, 128, 64)
1097
             self.final = nn.Conv2d(64, 4, 1)
1098
1099
         def forward(self, x):
1100
             A, H, V, D = self.dwt(x).chunk(4, 1)
1106
             a1 = self.eA1(A)
1102
1103
             d1 = self.eD1(torch.cat([H, V, D], 1))
1104
             d0 = torch.cat([a1, d1], 1)
1105
             p1 = self.pool(d0)
             pA, pD = p1[:, :64], p1[:, 64:]
1106
             a2 = self.eA2(pA)
1107
             d2 = self.eD2(pD)
1108
             d1b = torch.cat([a2, d2], 1)
1109
             b = self.b(self.pool(d1b))
1150
             u2 = self.u2(b, d1b)
1156
             u1 = self.u1(u2, d0)
11152
1113
             out = self.final(u1)
1154
             return self.idwt(out)
1115
11166
     # -- 4) OCRDataset for full training --
1167
11168
     class OCRDataset(Dataset):
11169
         def __init__(self, noisy, clean, pp, ps, mode):
1120
             self.noisy = noisy
             self.clean = clean
1126
             self.pp = pp
1122
1123
             self.ps = ps
             self.mode = mode
1124
             self.toT = transforms.ToTensor()
1125
1126
1127
         def __len__(self):
1128
             return len(self.noisy) * self.pp
1129
1130
         def __getitem__(self, idx):
1136
             ii = idx // self.pp
             n = Image.open(self.noisy[ii]).convert("L")
1132
```

```
c = Image.open(self.clean[ii]).convert("L")
1133
             tn, tc = self.toT(n), self.toT(c)
1134
             H, W = tn.shape[-2], tn.shape[-1]
1135
             y = random.randint(0, H - self.ps)
1136
1137
             x = random.randint(0, W - self.ps)
             tn = tn[:, y : y + self.ps, x : x + self.ps]
1138
             tc = tc[:, y : y + self.ps, x : x + self.ps]
1139
             if self.mode == "train":
1140
                 if random.random() < 0.5:</pre>
1146
1142
                     tn, tc = torch.flip(tn, [-1]), torch.flip(tc, [-1])
1143
                 if random.random() < 0.5:</pre>
                    tn, tc = torch.flip(tn, [-2]), torch.flip(tc, [-2])
1144
                 k = random.randint(0, 3)
1145
1146
                 if k > 0:
                    tn, tc = torch.rot90(tn, k, [1, 2]), torch.rot90(tc, k, [1, 2])
1147
1148
             return tn, tc
1149
1156
     # -- 5) Prepare train/val split --
1156
     all_no = sorted(glob.glob("./input/train/*.png"))
1152
     all_cl = sorted(glob.glob("./input/train_cleaned/*.png"))
1153
     idxs = list(range(len(all_no)))
1154
     random.seed(42)
1 1255
     random.shuffle(idxs)
     split = int(0.9 * len(idxs))
1257
     tr_idxs, va_idxs = idxs[:split], idxs[split:]
11258
     noisy_tr = [all_no[i] for i in tr_idxs]
11259
     clean_tr = [all_cl[i] for i in tr_idxs]
     noisy_va = [all_no[i] for i in va_idxs]
11266
     clean_va = [all_cl[i] for i in va_idxs]
11/62
11/63
     # -- 6) DataLoaders --
11964
     train_ds = OCRDataset(noisy_tr, clean_tr, pp=25, ps=256, mode="train")
     val_ds = OCRDataset(noisy_va, clean_va, pp=1, ps=256, mode="val")
1266
     train_loader = DataLoader(
1167
         train_ds, batch_size=8, shuffle=True, num_workers=4, pin_memory=True
1168
1169
1270
     val_loader = DataLoader(
         val_ds, batch_size=1, shuffle=False, num_workers=2, pin_memory=True
1276
1172
11273
     # -- 7) Model, optimizer, scheduler, loss w/ tuned weights --
     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
11275
     model = WaveletUNet().to(device)
11276
     dwt = HaarDWT().to(device)
11277
     sobel = Sobel().to(device)
     opt = torch.optim.AdamW(model.parameters(), lr=1e-3, weight_decay=1e-4)
11279
     sched = torch.optim.lr_scheduler.CosineAnnealingLR(opt, T_max=30)
11280
11286
     # Tuned weights
     alpha = 0.769634
11283
     beta = 0.430968
1184
     gamma = 0.250750
11/85
11286
11287
     def composite_loss(o, t):
11288
         11 = F.11_{loss(o, t)}
11/89
1196
         12 = F.mse_loss(dwt(o), dwt(t))
         e1 = F.l1_loss(sobel(o), sobel(t))
11296
1192
         return alpha * 11 + beta * 12 + gamma * e1
11293
1194
1195
     # -- 8) Train + validate --
11296
     best_rmse = 1e9
     patience = 5
1197
```

```
wait = 0
11298
     for epoch in range(1, 31):
1200
         model.train()
         tloss = 0.0
1201
         for xb, yb in train_loader:
1202
             xb, yb = xb.to(device), yb.to(device)
1203
1204
             pred = model(xb)
1205
             loss = composite_loss(pred, yb)
             opt.zero_grad()
1206
1207
             loss.backward()
1208
             opt.step()
             tloss += loss.item()
1209
         sched.step()
1210
1256
         # sliding-window validation
1212
1213
         model.eval()
         rmse_sum = 0.0
1254
         stride, ps = 128, 256
1215
         with torch.no_grad():
1216
1267
             for nf, cf in zip(noisy_va, clean_va):
                 imn = Image.open(nf).convert("L")
1218
                 imc = Image.open(cf).convert("L")
1219
                 tn = transforms.ToTensor()(imn).unsqueeze(0).to(device)
1220
                 tc = transforms.ToTensor()(imc).unsqueeze(0).to(device)
1221
1222
                 _{-}, _{-}, _{-}, _{-}, _{-} _{-} _{-}
1223
                 acc = torch.zeros_like(tn)
                 cnt = torch.zeros_like(tn)
1224
                 xs = list(range(0, W - ps + 1, stride)) + [W - ps]
1225
                 ys = list(range(0, H - ps + 1, stride)) + [H - ps]
1226
                 for x in xs:
1227
                     for y in ys:
1228
1229
                         p = tn[:, :, y : y + ps, x : x + ps]
                         out = model(p)
1230
1231
                         acc[:, :, y : y + ps, x : x + ps] += out
1232
                         cnt[:, :, y : y + ps, x : x + ps] += 1
1233
                 rec = acc / cnt
1234
                 mse = F.mse_loss(rec, tc).item()
1235
                 rmse_sum += math.sqrt(mse)
         val_rmse = rmse_sum / len(noisy_va)
1236
         print(
1237
             f"Epoch {epoch}: train_loss={tloss/len(train_loader):.4f}, val_RMSE={
1238
                  val_rmse:.6f}"
1239
1240
1241
         if val_rmse < best_rmse:</pre>
1242
             best_rmse = val_rmse
1243
1244
             wait = 0
             torch.save(model.state_dict(), "./working1/best.pth")
1245
         else:
1246
             wait += 1
1247
             if wait >= patience:
1248
                 print("Early stopping.")
1249
                 break
1250
1251
1252
     print("Best validation RMSE:", best_rmse)
1253
     # -- 9) Inference & submission --
1254
1255
     model.load_state_dict(torch.load("./working1/best.pth"))
1256
     test_files = sorted(glob.glob("./input/test/*.png"))
1257
     with open("./working1/submission.csv", "w") as fout:
1258
         fout.write("id, value\n")
1259
1260
         for tf in test_files:
1261
             im = Image.open(tf).convert("L")
             t = transforms.ToTensor()(im).unsqueeze(0).to(device)
1262
```

```
_, _, H, W = t.shape
1263
             acc = torch.zeros_like(t)
1264
             cnt = torch.zeros_like(t)
1265
             xs = list(range(0, W - ps + 1, stride)) + [W - ps]
1266
             ys = list(range(0, H - ps + 1, stride)) + [H - ps]
1267
             with torch.no_grad():
1268
                 for x in xs:
1269
                     for y in ys:
1270
                         p = t[:, :, y : y + ps, x : x + ps]
1271
1272
                         out = model(p)
1273
                        acc[:, :, y : y + ps, x : x + ps] += out
                         cnt[:, :, y : y + ps, x : x + ps] += 1
1274
             rec = (acc / cnt).clamp(0, 1).cpu().numpy()[0, 0]
1275
1276
             img_id = os.path.splitext(os.path.basename(tf))[0]
             for r in range(H):
1277
                 for c in range(W):
1278
                     fout.write(f''\{img_id\}_{r+1}_{c+1}, \{rec[r,c]:.6f\}\n'')
1279
1280
     print("Submission written to ./working1/submission.csv")
1282
```

C.5 Report Generation

1283

1284 The report generated by the second instance of CoMind is:

```
1285
      Our final pipeline is an Advanced Wavelet U-Net with CBAM attention modules, Composite Loss, and sliding-
1287
            window inference. Specifically:
1289
      1. Data Preparation
1290
         - Read all 112 noisy/clean train pairs, split 90/10 train/val by file indices (random seed 42).
1291
         - On-the-fly extract 25 random 256*256 patches per train image each epoch, with random horizontal/
1292
               vertical flips and 90 degree rotations.
1293
         - Validation uses a single full-image sliding-window pass (256 patches with 50% overlap) on the 10% hold
1294
               -out.
1295
1296
      2. Model Architecture
1297
         - 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)
1298
1299
               coefficients.
         - Two-level downsampling with MaxPool, interleaved CBAMResBlocks, and a 256-channel bottleneck.
1300
         - Decoder with ConvTranspose upsampling and CBAMResBlock skip-connections (UpRes), recombining wavelet
1301
1302
              bands.
1303
         - Final 1*1 convolution to 4 coefficient channels and inverse DWT back to image domain.
1304
1305
      3. Composite Loss
         - L1(pixels) weighted by alpha=0.769634
1306
           L2(wavelet coefficients) weighted by beta=0.430968
1307
1308
         - L1(edge Sobel) weighted by gamma=0.250750
1309
1310
      4. Training
         - Optimizer: AdamW(lr=1e-3, weight_decay=1e-4).
1311
1312
         - Scheduler: CosineAnnealingLR(T_max=30).
         - Batch size 8, 30 epochs, early stopping (patience=5) based on full-image RMSE on the val set.
1313
1314
1315
      5. Inference
         - Sliding-window full-image inference with 256*256 patches, stride=128, averaging overlapped pixels.
1316
1317
         - Generate submission.csv with id=image_row_col and float values.
1318
1319
      6. Hyperparameter tuning
1320
         - Optuna over (alpha, beta, gamma) on a 20-train/5-val mini-split for 10 trials, each with 2 epochs and
1321
              10 patches/image
1322
           Best found: alpha=0.769634, beta=0.430968, gamma=0.250750.
1323
      Final validation RMSE on 10% hold-out: ~0.012685.
1324
1325
1326
      References: CBAM (Woo et al., ECCV 2018), Haar DWT.
1327
1328
      Summary:
1329
       - Haar DWT:
1330
        Novelty: 2
1331
        Rationale: Uses standard 1-level Haar wavelet decomposition.
1332
        Rationale: Straightforward fixed conv implementation in PyTorch.
1333
1334
        Effectiveness: 6
```

```
1335
        Rationale: Explicit frequency separation aids denoising of high-frequency noise.
1336
        Efficiency: 9
1337
        Rationale: Lightweight fixed kernels on GPU.
1338
        Confidence: 8
        Rationale: Verified on toy and full runs.
1339
1340
1341
       - CBAMResBlock:
1342
        Novelty: 5
1343
        Rationale: Incorporates both channel & spatial attention for context-aware features.
        Feasibility: 8
1344
1345
        Rationale: Moderate complexity but supported by existing PyTorch modules.
1346
        Effectiveness: 7
1347
        Rationale: Improves focus on text strokes, empirically lowers RMSE.
1348
        Efficiency: 7
        Rationale: Adds overhead (~1.5*), but fits within GPU budget.
1349
1350
        Confidence: 9
        Rationale: Confirmed across toy and full training.
1351
1352
1353
      - UpRes (ConvTranspose + CBAM):
1354
        Novelty: 3
1355
        Rationale: Standard upsampling with skip-connections and attention.
1356
        Feasibility: 9
1357
        Rationale: Simple extension of ResBlock.
1358
        Effectiveness: 6
1359
        Rationale: Maintains detail in reconstruction.
1360
        Efficiency: 8
1361
        Rationale: Reasonable GPU use.
1362
        Confidence: 8
1363
        Rationale: Working in end-to-end script.
1364
1365
      - Wavelet UNet overall:
1366
        Novelty: 7
1367
        Rationale: Dual-branch U-Net in wavelet domain combining DWT, CBAM, skip connections.
1368
        Feasibility: 7
1369
        Rationale: Moderate engineering but fully implemented.
1370
        Effectiveness: 8
1371
        Rationale: Achieved RMSE ~0.0127, strong denoising.
1372
        Efficiency: 7
1373
        Rationale: 25 patches*8 batch; 30 epochs in ~20min.
1374
        Confidence: 8
        Rationale: Reproducible on A6000.
1375
1376
1377
      - Composite Loss:
1378
        Novelty: 6
1379
        Rationale: Balances pixel, frequency, and edge domains.
1380
        Feasibility: 8
        Rationale: Simple weighted sum.
1381
1382
        Effectiveness: 8
1383
        Rationale: Lower RMSE than L1 or L2 alone.
1384
        Efficiency: 7
        Rationale: Small overhead for extra transforms.
1385
1386
        Confidence: 9
1387
        Rationale: Verified in tuning and full runs.
1388
1389
      - Data Augmentation:
1390
        Novelty: 3
        Rationale: Flips and rotations.
1391
1392
        Feasibility: 10
        Rationale: Trivial to implement.
1393
1394
        Effectiveness: 5
1395
        Rationale: Standard but helpful for invariance.
1396
        Efficiency: 9
        Rationale: CPU-level only.
1397
        Confidence: 9
1398
1399
        Rationale: Common practice.
1400
1401
      - CosineAnnealingLR & EarlyStopping:
1402
        Novelty: 4
        Rationale: Standard learning-rate decay and val-based stop.
1403
1404
        Feasibility: 10
1405
        Rationale: Built-in PyTorch.
1406
        Effectiveness: 6
1407
        Rationale: Helps convergence and prevents overfit.
1408
        Efficiency: 9
1409
        Rationale: Minimal overhead.
1410
        Confidence: 9
1411
        Rationale: Observed smoothed training curves.
1412
1413 - Sliding-window Inference:
1414
        Novelty: 3
1415
        Rationale: Ensures full-image coverage.
```

```
1416
        Feasibility: 10
        Rationale: Straightforward loops.
1417
1418
        Effectiveness: 7
        Rationale: Accurate RMSE measurement, smooth outputs.
1419
1420
        Efficiency: 8
1421
        Rationale: Reasonable inference time (~10 min).
1422
        Confidence: 9
        Rationale: Verified end-to-end.
1423
1424
1425
       - Hyperparameter Tuning (Optuna):
1426
        Novelty: 4
1427
        Rationale: Automated search of loss weights.
1428
        Feasibility: 8
1429
        Rationale: Low-cost mini-training.
1430
        Effectiveness: 7
        Rationale: Found better weights vs. defaults.
1431
1432
        Efficiency: 5
1433
        Rationale: 10 trials * 2 epochs ~8 min.
1434
        Confidence: 7
1435
        Rationale: Verified chosen parameters.
1436
1437
      Weaknesses and Suggestions:
1438
      1. EMA & TTA Bug: The EMA + 4-way TTA implementation caused RMSE blow-ups, likely due to incorrect
1439
            accumulation. Before reintroducing, validate TTA on toy splits or isolate EMA without TTA.
1440
      2. Multi-Level Wavelets: Extend to 2-3 levels of DWT to capture coarser scales, with separate branches per
1441
            level.
1442
      3. Advanced Loss Terms: Incorporate perceptual (VGG) or structural (SSIM) losses to encourage visual
            fidelity.
1443
1444
      4. Dynamic Loss Weights: Schedule (alpha, beta, gamma) during training (e.g., start with high L1, gradually
1445
            emphasize wavelet).
1446
      5. Mixed Precision & Gradient Checkpointing: Further reduce memory and speed up with AMP and checkpointing
1447
1448
      6. Deeper/Ubiquitous Attention: Explore transformer-style attention (Swin, ViT) in encoder for global
1449
      7. Ensemble Models: Train multiple random seeds or architectures (ResUNet, UNet++) and average outputs.
1450
1451
      8. Multi-Scale Inference: Use patches at multiple scales or whole-image passes via down/upscaling for
1452
1453
      9. Adaptive Overlap: Use weights at edges of patches to smoothly blend overlaps (Gaussian windows).
      10. Distributed Hyperparameter Search: Scale Optuna across more trials or use BOHB for architecture search.
1454
1455
1456
      Implementing these will likely push RMSE even lower and improve robustness across document noise patterns.
1457
      Full code: ...(omitted)...
```

1458 C.6 The Second Iteration

Upon receiving two reports generated in the previous stage, CoMind proceeds to its second iteration. Following the brainstorming and idea refinement phases, the current idea pool is presented below. Notably, entries (3), (4), (5), and (7) are newly proposed by the agent, while the remaining ideas have been refined from those in the previous iteration.

- 1463 (0) Patch-based convolutional super-resolution denoiser: a sliding-window, patch-based model that predicts 1464 multiple denoised and super-resolved pixels per patch. Architecture inspired by Waifu2x with 1465 progressively increasing filter counts, LeakyReLU activations, and He initialization for robust 1466 convergence.
- 1467 (1) Advanced wavelet-domain dual-branch U-Net with CBAM attention and composite loss: decompose inputs via 1468 fixed 1-level DWT (LL, LH, HL, HH bands), encode approximation and detail separately with CBAM 1469 ResBlocks, decode and fuse via 1*1 convolutions, and train end-to-end using a weighted sum of pixel 1470 L1, wavelet-band L2, and edge L1 losses. Optimized with AdamW and cosine-annealing LR scheduling.
- 1471 (2) GAN-based restoration framework: a ResNet-based generator and 70*70 PatchGAN discriminator trained
 1472 with combined losses-L1 pixel loss, adversarial loss, stroke-consistency loss (via frozen stroke1473 feature CNN), and perceptual OCR-feature loss. Includes R1 gradient penalty and spectral
 1474 normalization for stability.
- 1475 (3) Masked autoencoder with vision transformer for denoising: patchify each image into non-overlapping
 1476 square tokens, randomly mask a high percentage, pretrain a ViT encoder (12 layers, hidden 768, 12
 1477 heads) plus light transformer decoder on L2 reconstruction of dirty images, then append an MLP head
 1478 and fine-tune end-to-end on noisy->clean pairs with L1 pixel + differentiable OCR-confidence loss.
 1479 Employ random block dropout and color jitter during fine-tuning; at inference use full-image encoding
 1480 or averaged mask schedules.
- (4) Conditional diffusion-based restoration: define a forward Gaussian-noise diffusion schedule, train a 5-1482 level U-Net conditioned on the dirty image via channel concatenation and FiLM/cross-attention of sinusoidal timestep embeddings. Use the standard DDPM MSE loss with classifier-free guidance, and sample with a deterministic DDIM sampler (~50 steps). Optionally post-process with bilateral or median filtering to remove speckles.
- 1486 (5) Learnable spectral gating in the Fourier domain: compute the 2D FFT of the dirty image, split its
 1487 spectrum into low/mid/high radial bands, apply learnable complex masks per band, and modulate each by
 1488 gate scalars predicted by a lightweight CNN on the dirty image. Recombine via inverse FFT and train
 1489 end-to-end with L2 pixel loss plus a spectral-smoothness regularizer on the masks.

```
1490 (6) Hypernetwork-modulated U-Net: extract per-image noise statistics (mean, std, skew, kurtosis, histogram
1491 bins), feed into an MLP hypernetwork that outputs FiLM scale (gamma) and shift (beta) parameters for
1492 selected convolutional feature maps of a base U-shaped CNN. Randomly augment noise levels during
1493 training; train end-to-end on noisy->clean with L1 loss and a small regularizer pushing gamma->1,
1494 beta->0. At inference compute stats per image, generate FiLM params, and denoise via the modulated U-
1495 Net.
```

1496

1497 1498

1499

1564

1565

1566

1567

/4.

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

- (7) Blind-spot test-time adaptation: for each test image, randomly mask a subset of pixels and fine-tune a lightweight CNN with skip connections on the single image using masked L2 loss and augmentations for a few hundred gradient steps. Optionally ensemble multiple random initializations to reduce variance.
- 1500 (8) Multi-model ensemble with diverse preprocessing: cluster trained models by their error patterns and
 1501 combine their outputs. Apply different preprocessing pipelines (raw, background-removed, edge maps,
 1502 thresholded) to the input, denoise with clustered sub-ensembles, and fuse predictions for robustness
 1503 across noise characteristics.
- 1504 (9) Enhanced augmentation and target refinement: simulate realistic scanning artifacts by applying 3D text
 1505 deformations, lighting variations, and shadows to clean images. Identify and compensate for
 1506 systematic artifacts in the provided 'clean' targets (e.g., single-pixel halos) by either modeling
 1507 them as noise or adjusting ground-truth labels during training.

```
And solution drafts generated in this iteration are:
1508
1509
      Pipeline 1: ResNet-34 Encoder U-Net with Multi-Scale Edge & Total-Variation Loss
1510
1511
1512
      A robust baseline using a pretrained ResNet-34 backbone as a U-Net encoder fused with a light-weight
1513
            decoder. Combines L1 loss, Sobel edge loss at multiple scales, and a total-variation regularizer to
1514
            preserve text strokes while smoothing background noise. Mixed precision training and sliding-window
1515
            inference ensure the entire pipeline runs in ~45 min on an A6000.
1516
1517
         Read all train noisy/clean PNGs, normalize to [0,1].
1518
       - Extract on-the-fly 256*256 patches: random crop + random horizontal/vertical flips + 90 degree rotations
1520
       - 90/10 split by file indices (seed=42). Use batch size 8-16.
1522
1523
      2. Model Architecture
1524
        - Encoder: torchvision.models.resnet34(pretrained=True), first conv modified to 1->64 channels.
1525
       - Decoder: four upsampling stages (ConvTranspose2d + Conv2d+BN+ReLU) mirroring ResNet blocks, with skip-
1526
            connections from encoder layers.
1527
       - Final conv 64->1 + Sigmoid.
1528
1529
      3. Loss Function
1530
       Let y_hat and y be predictions and targets.
1531
       - L1Loss(v_hat,v)
1532
       - Edge loss: L1 between Sobel(y_hat) and Sobel(y) at both full resolution and half resolution (downsample
1533
1534
       - TV: lambda*TV(v hat) where TV = mean(|\nabla xv hat|+|\nabla vv hat|).
       Total loss = alpha*L1 + beta*Edge_full + gamma*Edge_half + delta*TV, e.g. alpha=1.0, beta=0.5, gamma=0.25,
1535
             delta=1e-5.
1536
1537
1538
      4. Optimization
       - Optimizer: AdamW(lr=1e-3, weight_decay=1e-4).
1539
        - Scheduler: CosineAnnealingLR(T_max=25).
1540
1541
       - Mixed precision via torch.cuda.amp.
       - Early stopping on validation RMSE (patience=5).
1542
1543
1544
      5. Inference & Submission
1545
       - Perform sliding-window inference on each test image with 256*256 patches, stride=128.
1546
       - Average overlapping patches.
1547
       - Clamp outputs to [0,1], write submission.csv with id=image_row_col.
1548
      Compute budget: ~20 min train + ~5 min inference.
1549
1550
1551
      Pipeline 2: Laplacian-Pyramid Multi-Scale Residual U-Net with Pyramid Loss
1552
1553
1554
      A novel pyramid-domain network that decomposes images into multi-scale Laplacian bands, denoises each band
1555
            via shared-weight residual blocks, and merges them back. Multi-level L1 losses focus the model on
1556
            both coarse structures and fine text details. Efficient and fully end-to-end in PyTorch, training
            finishes in ~50 min on an A6000.
1557
1558
1559
      1. Data & Augmentation
1560
       - Same data split and patch sampling (256*256, flips, rotations).
1561
      2. Laplacian Pyramid Transform (LPT)
1562
1563
       - On the fly, for each patch: create 2-level Gaussian pyramid using avg-pooling (scale 1->0.5->0.25), then
```

 $\label{eq:compute Laplacian bands L0=(I-upsample(I/2)), L1=(I/2-upsample(I/4)), and a residual low band L2=II and a residual low band l$

```
1568
      3. Network
1569
       - Encoder-decoder U-Net with 3-channel input and 3-channel output, 4 down/upsampling stages,
             64->128->256->512->256->128->64 channel widths.
1570
       - Each block: two 3*3 Conv+BN+ReLU and a residual skip connection.
1571
1572
1573
      4. Reconstruction
       - Output is 3 residual corrections RO,R1,R2 to each band.
1574
       - Reconstruct clean patch:
1575
         a0 = L0 + R0; a1 = L1 + R1; a2 = L2 + R2;
I = a0 + up(a1) + up(up(a2)).
1576
1577
1578
1579
      5. Loss Function
1580
       - Per-band L1 losses: w0*L1(a0,H0) + w1*L1(a1,H1) + w2*L1(a2,H2), where Hk are ground-truth bands and w0
            =1.0, w1=0.5, w2=0.25.
1581
       - Pixel L1: lambda_pix*L1(I, Igt) with lambda_pix=0.1.
1582
1583
      6. Training
1584
       - AdamW(lr=5e-4, weight_decay=1e-5), one-cycle scheduler (max_lr=5e-4, epochs=30).
1585
1586
       - Batch size 8, mixed precision.
1587
       - Early stopping on val RMSE (patience=6).
1588
1589
      7. Inference
       - On full test images, run sliding-window with 256*256 patches, stride=128.
1590
       - For each patch, apply pyramid -> model -> reconstruct -> accumulate and average.
1591
1592
       - Clamp to [0,1], write submission.csv.
1593
      Compute budget: ~25 min training + ~8 min inference.
1594
```

The best code proposed in the second iteration is:

```
1595
     import os, random, csv
1597
     import numpy as np
1598
     from glob import glob
1599
1600
     from PIL import Image
1601
     import torch
1602
     import torch.nn as nn
1603
     import torch.nn.functional as F
1604
     from torch.utils.data import Dataset, DataLoader
1605
     from torchvision import transforms, models as tv_models
1606
     from torch.cuda.amp import autocast, GradScaler
1607
1608
     # 1) Reproducibility
1609
     SEED = 42
1610
     random.seed(SEED)
1611
     np.random.seed(SEED)
1612
     torch.manual_seed(SEED)
     torch.cuda.manual_seed_all(SEED)
1614
1615
     DEVICE = torch.device("cuda" if torch.cuda.is_available() else "cpu")
1616
1617
1618
     # 2) Paths
     TRAIN_NOISY = "./input/train"
1619
     TRAIN_CLEAN = "./input/train_cleaned"
1620
     TEST_DIR = "./input/test"
1621
     WORK_DIR = "./working0"
1622
1623
     os.makedirs(WORK_DIR, exist_ok=True)
1624
     # 3) Hyperparameters
1625
1626
     PATCH_SIZE = 256
     STRIDE = 128
1627
     BATCH_SIZE = 12
1628
1629
     LR, WD = 1e-3, 1e-4
     MAX\_EPOCHS = 50
1630
     PATIENCE = 10
1631
1632
     T_MAX = 50 \# for LR scheduler
1633
     # 4) Loss-weight constants (from tuning)
1634
     w1, w2, w3, w4 = 1.0, 0.5, 0.25, 1e-5
1635
    lambda_aux = 0.4394633936788146
1636
```

```
lambda_mse = 0.1312037280884873
1637
     lambda_ssim = 0.031198904067240532
     lambda_ssim2 = lambda_ssim / 2
1639
1640
1641
     # 5) Dataset + augmentations
1642
     class OCRDataset(Dataset):
1643
         def __init__(self, noisy_list, clean_list, ps, train):
1644
             self.noisy, self.clean = noisy_list, clean_list
1645
1646
             self.ps, self.train = ps, train
1647
             self.to_tensor = transforms.ToTensor()
1648
             self.aug = transforms.Compose(
                 Γ
1649
                     transforms.RandomChoice(
1650
                         Γ
1651
                             transforms.RandomHorizontalFlip(1.0),
1652
                             transforms.RandomVerticalFlip(1.0),
1653
                             transforms.RandomRotation(90),
1654
                             transforms.RandomRotation(180),
1655
1656
                             transforms.RandomRotation(270),
                         1
1657
                     ),
1658
                     transforms.RandomApply([transforms.GaussianBlur(3, (0.1, 2.0))], p
1659
1660
                     transforms.RandomApply([transforms.RandomAdjustSharpness(2.0)], p
1661
1662
                          =0.3).
                 ]
1663
             )
1664
1665
         def __len__(self):
1666
             return len(self.noisy)
1667
1668
         def __getitem__(self, i):
1669
             n = Image.open(self.noisy[i]).convert("L")
1670
             c = Image.open(self.clean[i]).convert("L")
1671
1672
             w, h = n.size
1673
1674
             if w < self.ps or h < self.ps:</pre>
                 pad = (0, 0, \max(0, \text{self.ps - w}), \max(0, \text{self.ps - h}))
1675
                 n = transforms.functional.pad(n, pad, fill=255)
1676
                 c = transforms.functional.pad(c, pad, fill=255)
1677
                 w, h = n.size
1678
             # crop
1679
             if self.train:
1680
                 x = random.randint(0, w - self.ps)
1681
1682
                 y = random.randint(0, h - self.ps)
1683
             else:
                 x = (w - self.ps) // 2
1684
                 y = (h - self.ps) // 2
1685
             n = n.crop((x, y, x + self.ps, y + self.ps))
1686
             c = c.crop((x, y, x + self.ps, y + self.ps))
1687
             if self.train and random.random() < 0.5:</pre>
1688
                 n = self.aug(n)
1689
                 c = self.aug(c)
1690
1691
             return self.to_tensor(n), self.to_tensor(c)
1692
1693
1694
      # 6) Prepare train/val split
     noisy_files = sorted(glob(f"{TRAIN_NOISY}/*.png"))
1695
     {\tt clean\_files = [f"{TRAIN\_CLEAN}/" + os.path.basename(x) \ for \ x \ in \ noisy\_files]}
1696
1697
     N = len(noisy_files)
     idx = list(range(N))
1698
1699
     random.shuffle(idx)
1700
     ntr = int(0.9 * N)
    tr_idx, va_idx = idx[:ntr], idx[ntr:]
1703
```

```
train_noisy = [noisy_files[i] for i in tr_idx]
1702
     train_clean = [clean_files[i] for i in tr_idx]
     val_noisy = [noisy_files[i] for i in va_idx]
1704
     val_clean = [clean_files[i] for i in va_idx]
1705
1706
     train_ds = OCRDataset(train_noisy, train_clean, PATCH_SIZE, train=True)
1707
1708
     val_ds = OCRDataset(val_noisy, val_clean, PATCH_SIZE, train=False)
     train_loader = DataLoader(
1709
         train_ds, batch_size=BATCH_SIZE, shuffle=True, num_workers=4, pin_memory=True
1710
1713
1712
     val_loader = DataLoader(
         val_ds, batch_size=BATCH_SIZE, shuffle=False, num_workers=4, pin_memory=True
1713
1714
1715
     # 7) Sobel, TV, SSIM helpers
1716
1717
     sob_x = (
         torch.tensor([[1, 0, -1], [2, 0, -2], [1, 0, -1]], dtype=torch.float32)
1718
         .view(1, 1, 3, 3)
1719
         .to(DEVICE)
1720
1723
     sob_y = sob_x.transpose(2, 3)
1722
1723
1724
     def sobel(x):
1725
         gx = F.conv2d(x, sob_x, padding=1)
1726
         gy = F.conv2d(x, sob_y, padding=1)
1727
         return torch.sqrt(gx * gx + gy * gy + 1e-6)
1728
1729
1730
     def total_variation(x):
1731
         dh = (x[:, :, 1:, :] - x[:, :, :-1, :]).abs().mean()
1732
         dw = (x[:, :, :, 1:] - x[:, :, :, :-1]).abs().mean()
1733
         return dh + dw
1734
1735
1736
     def ssim_map(a, b, C1=0.01**2, C2=0.03**2):
1737
1738
         mu_a = F.avg_pool2d(a, 3, 1, 1)
1739
         mu_b = F.avg_pool2d(b, 3, 1, 1)
         sa = F.avg_pool2d(a * a, 3, 1, 1) - mu_a * mu_a
1740
         sb = F.avg_pool2d(b * b, 3, 1, 1) - mu_b * mu_b
1741
         sab = F.avg_pool2d(a * b, 3, 1, 1) - mu_a * mu_b
1742
1743
         num = (2 * mu_a * mu_b + C1) * (2 * sab + C2)
         den = (mu_a * mu_a + mu_b * mu_b + C1) * (sa + sb + C2)
1744
         return num / (den + 1e-8)
1745
1746
1747
     def ssim_loss(a, b):
1748
         return 1.0 - ssim_map(a, b).mean()
1749
1750
1753
     # 8) loss_terms
1752
     11_loss = nn.L1Loss()
1753
     mse_loss = nn.MSELoss()
1754
1755
1756
     def loss_terms(pred, target):
1757
         L1v = 11_loss(pred, target)
1758
1759
         MSEv = mse_loss(pred, target)
         Ef = 11_loss(sobel(pred), sobel(target))
1760
         p2, t2 = F.avg_pool2d(pred, 2), F.avg_pool2d(target, 2)
1761
1762
         Eh = 11_{loss(sobel(p2), sobel(t2))}
         TVv = total_variation(pred)
1763
1764
         return L1v, MSEv, Ef, Eh, TVv
1765
1766
```

```
# 9) Model w/ deep supervision
1767
     class ResUNetDS(nn.Module):
1768
1769
         def __init__(self):
             super().__init__()
1770
1771
             r34 = tv_models.resnet34(pretrained=True)
             self.enc0 = nn.Conv2d(1, 64, 7, 2, 3, bias=False)
1772
             self.enc0.weight.data = r34.conv1.weight.data.mean(dim=1, keepdim=True)
1773
             self.bn0, self.relu0, self.pool0 = r34.bn1, r34.relu, r34.maxpool
1774
             self.enc1, self.enc2 = r34.layer1, r34.layer2
1775
1776
             self.enc3, self.enc4 = r34.layer3, r34.layer4
1777
1778
             def up(i, o):
                 return nn.ConvTranspose2d(i, o, 2, 2)
1779
1780
             def cb(i, o):
1783
1782
                 return nn.Sequential(
                    nn.Conv2d(i, o, 3, 1, 1, bias=False),
1783
                    nn.BatchNorm2d(o),
1784
                     nn.ReLU(inplace=True),
1785
                     nn.Conv2d(o, o, 3, 1, 1, bias=False),
1786
                    nn.BatchNorm2d(o),
1787
                    nn.ReLU(inplace=True),
1788
                 )
1789
1790
             self.up4, self.dec4 = up(512, 256), cb(256 + 256, 256)
1791
             self.up3, self.dec3 = up(256, 128), cb(128 + 128, 128)
1792
             self.up2, self.dec2 = up(128, 64), cb(64 + 64, 64)
1793
             self.aux\_up, self.aux\_out = up(64, 64), nn.Conv2d(64, 1, 1)
1794
1795
             self.up1, self.dec1 = up(64, 64), cb(64 + 64, 64)
             self.up0, self.outc = up(64, 64), nn.Conv2d(64, 1, 1)
1796
             self.sig = nn.Sigmoid()
1797
1798
         def forward(self, x):
1799
             x0 = self.relu0(self.bn0(self.enc0(x)))
1800
             x1 = self.pool0(x0)
1801
1802
             x2 = self.enc1(x1)
1803
             x3 = self.enc2(x2)
1804
             x4 = self.enc3(x3)
             x5 = self.enc4(x4)
1805
1806
             d4 = self.dec4(torch.cat([self.up4(x5), x4], dim=1))
1807
             d3 = self.dec3(torch.cat([self.up3(d4), x3], dim=1))
1808
             d2 = self.dec2(torch.cat([self.up2(d3), x2], dim=1))
1809
             aux = self.sig(self.aux_out(self.aux_up(d2)))
1810
             d1 = self.dec1(torch.cat([self.up1(d2), x0], dim=1))
1813
1812
             main = self.sig(self.outc(self.up0(d1)))
1813
             return main, aux
1814
1815
     model = ResUNetDS().to(DEVICE)
1816
1817
1818
     # 10) Optimizer, scheduler, scaler
     optimizer = torch.optim.AdamW(model.parameters(), lr=LR, weight_decay=WD)
1819
     scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, T_max=T_MAX)
1820
1821
     scaler = GradScaler()
1822
     # 11) Training + snapshot saving
1823
1824
     best_rmse = float("inf")
     patience = 0
1825
     snap\_epochs = set([10, 20, 30, 40, 50])
1826
1827
     for epoch in range(1, MAX_EPOCHS + 1):
1828
1829
         model.train()
1830
         train_loss = 0.0
1831
         for noisy_img, clean_img in train_loader:
```

```
noisy_img, clean_img = noisy_img.to(DEVICE), clean_img.to(DEVICE)
1832
1833
             optimizer.zero_grad()
             with autocast():
1834
                 main_pred, aux_pred = model(noisy_img)
1835
                 L1v, MSEv, Ef, Eh, TVv = loss_terms(main_pred, clean_img)
1836
                 s1 = ssim_loss(main_pred, clean_img)
1837
1838
                 p2, t2 = F.avg_pool2d(main_pred, 2), F.avg_pool2d(clean_img, 2)
                 s2 = ssim_loss(p2, t2)
1839
                 main_loss = (
1840
                     w1 * L1v
1841
1842
                     + lambda_mse * MSEv
                     + w2 * Ef
1843
                     + w3 * Eh
1844
                     + w4 * TVv
1845
                     + lambda_ssim * s1
1846
1847
                     + lambda_ssim2 * s2
                 )
1848
                 aux_up = F.interpolate(
1849
                     aux_pred,
1850
1851
                     size=clean_img.shape[-2:],
                     mode="bilinear",
1852
                     align_corners=False,
1853
1854
                 La, Ma, Ea, Eh2, TVa = loss_terms(aux_up, clean_img)
1855
1856
                 sa = ssim_loss(aux_up, clean_img)
                 pa, ca = F.avg_pool2d(aux_up, 2), F.avg_pool2d(clean_img, 2)
1857
                 sa2 = ssim_loss(pa, ca)
1858
                 aux_loss = (
1859
                     w1 * La
1860
                     + lambda_mse * Ma
1861
                     + w2 * Ea
1862
                     + w3 * Eh2
1863
                     + w4 * TVa
1864
1865
                     + lambda_ssim * sa
                     + lambda_ssim2 * sa2
1866
1867
1868
                 loss = main_loss + lambda_aux * aux_loss
1869
             scaler.scale(loss).backward()
             scaler.step(optimizer)
1870
             scaler.update()
1871
             train_loss += loss.item()
1872
1873
         scheduler.step()
1874
         # validation
1875
         model.eval()
1876
1877
         se, count = 0.0, 0
1878
         with torch.no_grad():
             for noisy_img, clean_img in val_loader:
1879
                 noisy_img, clean_img = noisy_img.to(DEVICE), clean_img.to(DEVICE)
1880
1881
                 with autocast():
                     pred, _ = model(noisy_img)
1882
                 se += ((pred - clean_img) ** 2).sum().item()
1883
                 count += pred.numel()
1884
         val_rmse = np.sqrt(se / count)
1885
1886
         print(
             f"Epoch {epoch}: TrainLoss={train_loss/len(train_loader):.4f}, ValRMSE={
1887
                  val_rmse:.6f}"
1888
1889
1890
1891
         # best + snapshot
1892
         if val_rmse < best_rmse:</pre>
             best_rmse = val_rmse
1893
1894
             torch.save(model.state_dict(), os.path.join(WORK_DIR, "best_full.pth"))
             patience = 0
1895
1896
         else:
```

```
patience += 1
1897
         if epoch in snap_epochs:
1898
             torch.save(model.state_dict(), os.path.join(WORK_DIR, f"snap_{epoch}.pth"))
1899
         if patience >= PATIENCE:
1900
             print("Early stopping.")
1901
1902
             break
1903
     print("Best validation RMSE:", best_rmse)
1904
1905
     # 12) Ensemble load
1906
1907
     ckpts = ["best_full.pth"] + sorted(
         [f for f in os.listdir(WORK_DIR) if f.startswith("snap_")],
1908
         key=lambda x: int(x.split("_")[1].split(".")[0]),
1909
1910
     )[-2:]
     ensemble_nets = []
1911
1912
     for ck in ckpts:
         net = ResUNetDS().to(DEVICE)
1913
         net.load_state_dict(torch.load(os.path.join(WORK_DIR, ck)))
1914
1915
1916
         ensemble_nets.append(net)
1917
1918
1919
     # 13) Sliding-window ensemble inference
     def ensemble_infer(img_arr):
1920
1921
         h, w = img_arr.shape
         inp = torch.from_numpy(img_arr / 255.0).unsqueeze(0).unsqueeze(0).to(DEVICE)
1922
         ph = (PATCH_SIZE - h % STRIDE) % STRIDE
1923
         pw = (PATCH_SIZE - w % STRIDE) % STRIDE
1924
         inp = F.pad(inp, (0, pw, 0, ph), mode="reflect")
1925
         _, _, H, W = inp.shape
1926
         out = torch.zeros_like(inp)
1927
         wt = torch.zeros_like(inp)
1928
         for y in range(0, H - PATCH_SIZE + 1, STRIDE):
1929
             for x in range(0, W - PATCH_SIZE + 1, STRIDE):
1930
                 patch = inp[:, :, y : y + PATCH_SIZE, x : x + PATCH_SIZE]
1931
                 preds = []
1932
1933
                 with torch.no_grad(), autocast():
1934
                     for net in ensemble_nets:
                        p, _ = net(patch)
1935
                        preds.append(p)
1936
                 avg_p = torch.stack(preds, 0).mean(0)
1937
                 out[:, :, y : y + PATCH_SIZE, x : x + PATCH_SIZE] += avg_p
1938
                 wt[:, :, y : y + PATCH\_SIZE, x : x + PATCH\_SIZE] += 1.0
1939
         out = out / wt
1940
         out = out[:, :, :h, :w]
1941
1942
         return out.detach().cpu().numpy().squeeze()
1943
1944
     # 14) Write submission.csv
1945
     submission_path = os.path.join(WORK_DIR, "submission.csv")
1946
     with open(submission_path, "w", newline="") as f:
1947
1948
         writer = csv.writer(f)
         writer.writerow(["id", "value"])
1949
         for tf in sorted(
1950
1951
             glob(f"{TEST_DIR}/*.png"), key=lambda x: int(os.path.basename(x).split(".")
1952
                  [0]
         ):
1953
1954
             img_id = os.path.basename(tf).split(".")[0]
             img = np.array(Image.open(tf).convert("L"), dtype=np.float32)
1955
1956
             den = ensemble_infer(img)
1957
             H, W = den.shape
             for i in range(H):
1958
1959
                 for j in range(W):
                     writer.writerow([f"{img_id}_{i+1}_{j+1}", f"{den[i,j]:.6f}"])
1960
     print("Submission saved to", submission_path)
```