

Memory-Augmented Agent Training for Business Document Understanding

Jiale Liu¹, Yifan Zeng², Malte Højmark-Bertelsen³, Marie Normann Gadeberg³,
Huazheng Wang², Qingyun Wu¹

¹Pennsylvania State University ²Oregon State University ³Beyond Work

Abstract

Traditional enterprises face significant challenges in processing business documents, where tasks like extracting transport references from invoices remain largely manual despite their crucial role in logistics operations. While Large Language Models offer potential automation, their direct application to specialized business domains often yields unsatisfactory results. We introduce Matrix (Memory-Augmented agent Training through Reasoning and Iterative eXploration), a novel paradigm that enables LLM agents to progressively build domain expertise through experience-driven memory refinement and iterative learning. To validate this approach, we collaborate with one of the world’s largest logistics companies to create a dataset of Universal Business Language format invoice documents, focusing on the task of transport reference extraction¹. Experiments demonstrate that Matrix outperforms prompting a single LLM by 30.3%, vanilla LLM agent by 35.2%. We further analyze the metrics of the optimized systems and observe that the agent system requires less API calls, fewer costs and can analyze longer documents on average. Our methods establish a new approach to transform general-purpose LLMs into specialized business tools through systematic memory enhancement in document processing tasks.

1 Introduction

Combing through large quantities of unstructured data remains a widespread challenge in enterprise operations, particularly in finance functions where efficient invoice processing represents a growing competitive advantage. Despite the prevalent adoption of digital invoicing, many organizations still grapple with the labor-intensive and error-prone task of manually extracting crucial identifiers from business transactions. For logistics companies manual extraction not only slows down operations but also introduces the potential for human error, leading to misrouted shipments and customer dissatisfaction.

Large Language Models have demonstrated remarkable capabilities in natural language understanding and processing (Achiam et al. 2023; Yang et al. 2024; Dubey et al.

2024). Building upon these capabilities, agent-based systems that leverage LLMs for multi-step reasoning have emerged as a promising paradigm (Hu, Lu, and Clune 2024; Wu et al. 2024a, 2023). However, these systems face significant challenges when adapting to specialized contexts without manual tuning, as LLMs are not specifically trained to reason over business documents.

To address these limitations, we propose Matrix (Memory-Augmented agent Training through Reasoning and Iterative eXploration), a novel framework that enables LLM-based agents to learn and adapt from their experiences. Matrix incorporates a unique iterative self-refinement mechanism that allows agents to systematically improve their understanding of document structures and extraction patterns. The agents go through task exploration and optimization to iteratively improve their insights on the general task structure. The distilled actionable insights will be saved into a long-term memory of the agents and will guide future task solving attempts, leading to agent systems with better strategic reasoning and task solving capability.

In collaboration with Kuehne+Nagel, one of the world’s leading logistics companies, we collect data from real-world business invoice documents and benchmark Matrix on the dataset. Specifically, we focus on the challenge of extracting transport reference numbers, which directly impact shipment routing and supply chain visibility. To facilitate broader research in this domain, we also introduce the first open-source dataset of anonymized business invoices. This benchmark poses a unique opportunity to evaluate agent system’s adaptive learning capabilities in critical business context.

Our experimental results demonstrate Matrix’s effectiveness after several epochs of optimization. The system achieves significant improvements through iterative learning and demonstrates strong performance. When compared to previous baselines, Matrix achieves 30.3% performance improvement compared to chain-of-thought prompting (Wei et al. 2022), surpasses vanilla LLM agents (Wu et al. 2023) by 35.2%, reflexion (Shinn et al. 2024) by 27.28%. Further analysis on the system reveals that the learning process leads to significantly lower latency and costs, with a stronger capability in processing longer documents.

The key contributions of this work are threefold:

1. We propose a novel agent-based system for document

¹The dataset contains sensitive information from corporate customers and therefore cannot be released. We release an anonymized subset to facilitate research in this field. Data available at <https://tinyurl.com/3pak9t4>

reasoning and information extraction. It serves as a strong baseline in business cases.

2. We introduce Matrix, a novel paradigm that enables LLM-based agents to systematically learn and adapt to novel tasks through iterative self-refinement and experience-based memory updating.
3. We present the first open-source benchmark for evaluating document understanding systems on standardized business documents, facilitating reproducible research and practical applications in enterprise settings. This benchmark provides a challenging and practically relevant test for document reasoning and information retrieval capabilities in business contexts using agentic networks.

2 Related Work

Business Document Reasoning Benchmarks. While general-domain question answering datasets have driven advances in natural language understanding (Abujabal et al. 2019; Rajpurkar et al. 2016; Yang et al. 2018; Joshi et al. 2017; Talmor and Berant 2018; Bajaj et al. 2016; Kwiatkowski et al. 2019), they fail to capture the structured, transactional nature of business documents. Existing research on business document understanding mostly focus on vision-based information extraction from scanned copies (Harley, Ufkes, and Derpanis 2015; Riba et al. 2019; Zhong, Tang, and Yepes 2019; Antonacopoulos et al. 2009), which emphasize layout understanding over semantic reasoning. These datasets, while valuable for OCR and structure recognition tasks, fail to capture the domain-specific patterns essential for business document processing. Moreover, text-based researches tend to use proprietary, undisclosed datasets (Hamdi et al. 2021; Krieger et al. 2021; Palm, Winther, and Laws 2017; Tarawneh et al. 2019), limiting reproducibility and real-world applicability. To address these limitations, we present a publicly available business invoice dataset derived from real-world transactions to facilitate research on practical business scenarios.

Prompt Optimization. Prompt optimization is a popular paradigm for maximizing LLM’s performance to novel tasks without expensive model tuning by finding a optimal task prompt (Zhou et al. 2022b; Pryzant et al. 2023; Cheng et al. 2023; Prasad et al. 2022). In-context learning emerges as a prominent paradigm, where a set of input-output pairs is provided as few shot examples to the LLM (Min et al. 2021; Dong et al. 2022; Brown 2020). By automatically retrieving demonstrations from training set (Zhao et al. 2021; Lu et al. 2021; Liu et al. 2021) or from adaptively annotated samples (Zhang et al. 2023b; Wu et al. 2022; Su et al. 2022), they primarily treat demonstrations as static examples rather than distilled insights. In contrast to this line of work, our proposed paradigm aims to distill trajectories into generalizable heuristics.

Agent Learning. There has been efforts in exploring inference time performance boost since the emergence of Large Language Models (Shinn et al. 2024; Madaan et al. 2024; Yao et al. 2023, 2024; Sumers et al. 2023; Wei et al. 2022; Zhou et al. 2022a; Guo et al. 2024). Recent works

have extended this paradigm to agentic systems. Some works represent and learn the optimal workflow of agentic systems in the form of complex graphs (Zhuge et al. 2024; Wu et al. 2024b), code (Hu, Lu, and Clune 2024), memory (Wang et al. 2024) and trees (Zhang et al. 2024a) to improve the system’s performance on complex tasks, while others learns reusable tools (Zhang et al. 2024c; Cai et al. 2023; Qian et al. 2023; Yuan et al. 2023) and experience (Zhao et al. 2024; Wang et al. 2024) for agentic systems. Different from previous memory optimization based works, our proposed approach does not rely on any pre-defined in-context examples (Zhao et al. 2024) and meanwhile guarantee the diversity and robustness of the learned memory.

3 Matrix: Memory-Augmented agent Training through Reasoning and Iterative eXploration

In this section, we begin with presenting the formulation of document reasoning, then introduce the framework of Matrix.

3.1 Problem Statement

We begin with the formal definition of the document reasoning problem. In this problem, we assume the existence of a dataset for a given task, which consists of N instances of document, query, and answer triplets:

$$\mathcal{D} = \{(d_i, q_i, a_i)\}_{i=1}^N \quad (1)$$

where d_i represents a document, q_i is an associated query, and a_i is the corresponding correct answer. In the context of transport reference extraction from business documents, this formulation naturally maps to our task of interest: each document d_i represents a UBL invoice, the query q_i requests the location of the transport reference number, and the answer a_i is the correct reference number. This mapping allows us to frame the challenge of invoice processing as a structured document reasoning problem while maintaining the generality of our approach.

The overall goal in document reasoning is to train a system that learns a function $f : D \times Q \rightarrow A$, where D, Q, A represents the space of documents, queries, and answers respectively. For a given document d_i and query q_i , the agent generates an answer \hat{a}_i . The system’s performance is evaluated based on the alignment between \hat{a}_i and a_i .

We consider an agent with a long-term memory module M , which stores useful contextual information for reasoning. For a given query (d, q) , the agent solves the task in an iterative manner. At each timestep t , the agent observes the current state o_t , takes an action a_t (which may involve interacting with the document or recalling information from memory), and then receives an updated observation o_{t+1} . This sequence of interactions produces a trajectory of observations and actions: $\tau = (o_0, a_0, o_1, a_1, \dots, o_T, a_T)$. Finally, the agent produces the final predicted answer \hat{a}_i through an answer extraction module $g(\tau)$, which takes into account the full trajectory τ of observations and actions

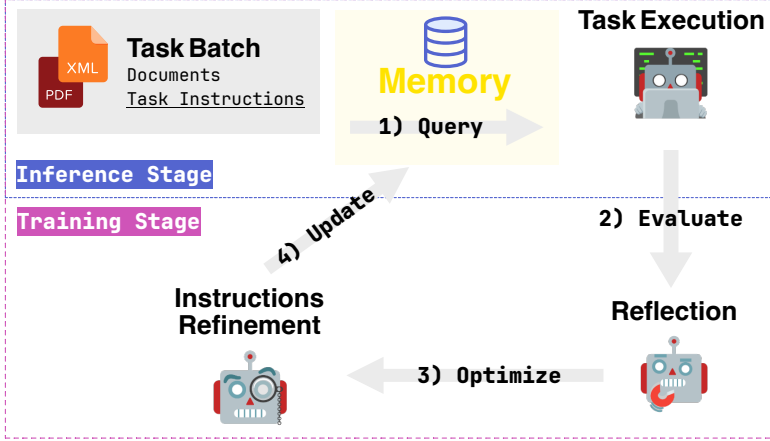


Figure 1: The training and inference pipeline of Matrix.

across timesteps. Thus, the system can be viewed as generating the prediction:

$$\hat{a}_i = g(\tau) = f(d_i, q_i | M) \quad (2)$$

The formal goal of optimizing the agent’s reasoning ability can be described via the objective:

$$\arg \min_M \mathbf{E}_{(d_i, q_i, a_i) \sim \mathcal{D}} [\mathcal{L}(f(d_i, q_i | M), a_i)] \quad (3)$$

where \mathcal{L} measures the discrepancy between generated answer and the ground truth. In our case, \mathcal{L} is defined using an exact match comparison.

3.2 Proposed Method

To optimize the objective in Eq.3, we follow the paradigm of supervised learning and develop the dataset into train set \mathcal{D}_{train} and test set \mathcal{D}_{test} to approximate the inaccessible real data distribution. We train an optimal memory and test it on test task set. The pipeline is illustrated in Figure 1.

Trajectory Sampling. We perform optimization by progressively updating the memory module M over multiple epochs. During each epoch, we sample a mini-batch of tasks from \mathcal{D}_{train} . For each sampled task (d_i, q_i, a_i) , the agent interacts through a sequence of actions over time, forming a trajectory $\tau_i = (o_0^i, a_0^i, o_1^i, a_1^i, \dots, o_T^i, a_T^i)$, where o_t^i and a_t^i denote the observation and the action at timestep t for the i -th task.

The trajectory continues until the agent either reaches a solution or the interaction exceeds a predefined maximum number of steps. We enforce a upper limit T_{max} to the total number of steps per task. This can constrain the token count in each trajectory, and allow us to increase the mini-batch size during optimization. By doing so, the optimizer can process more diverse tasks within a single batch, which increases the model’s exposure to varied problem domains, leading to a more robust and generalizable memory M .

Reflection. Despite the impressive reasoning capabilities of LLM-based agents, they are prone to issues such as hallucination (Ji et al. 2023), factual errors (Wang et al. 2023),

and reasoning failures (Huang et al. 2023). In a limited number of steps, it is hard for agent to autonomously recognize and rectify its own mistakes without explicit self-correction mechanisms or auxiliary prompts (Huang et al. 2023; Jiang et al. 2024; Pan et al. 2023; Kamoi et al. 2024; Song et al. 2024). To overcome this limitation, we introduce a Reflector module that operates as a post hoc evaluator of the trajectory. The reflector is provided with the trajectory τ_i and the ground truth answer y_i . Its task is to label the solution as ”Correct/Incorrect” and identify the key steps where reasoning errors occurred or correct decisions were made that led to the proper solution. The reflection process is described as:

$$r_i = \text{LM}_{\text{reflect}}(\tau_i, y_i) \quad (4)$$

where $\text{LM}_{\text{reflect}}(\cdot)$ represents prompting a LLM to evaluate the trajectory. The reflection phase provides insights for refining the memory and improving the agent’s task-solving strategy.

Optimization. To obtain an optimal solution in Eq.3, an optimizer is needed that can generate new solutions based on performance measurement. The optimization problem operates on the space of natural languages, which perfectly paves the way for borrowing Large Language Models’ exceptional capability in natural language understanding (Yang et al. 2023; Zhang et al. 2024c). We propose a meta-optimizer with Large Language Models as the backend. The optimizer takes in three parts of information: (1) The execution trajectory of the agents in solving the tasks in the current training mini-batch (2) The assessment of each trajectory provided by the evaluator. (3) The current memory the agents operate on. Although the memory can be accessed by reading the prompt in the trajectory, we explicitly provide the current memory to let the optimizer progressively update the instruction based on the assessments for improvement. The optimization process can be formulated as:

$$M_{i+1} = \text{LM}_{\text{optim}}(\tau, y, M_i) \quad (5)$$

where $\text{LM}_{\text{optim}}(\cdot)$ denotes the process of LLM-based optimization.

4 Evaluations

4.1 Experiment Setup

Agent system. We experiment with refining a two-agent system. In the system, an assistant powered by LLM receives and analyzes the task and suggests code for execution, and a user proxy automatically executes code and provide the result. The conversation will end once a final answer is reached by assistant or the max number of conversation turns is reached. The system is implemented by AG2² (formerly AutoGen) (Wu et al. 2023).

Experiment Dataset. We evaluate our system on a collection of real-world Universal Business Language invoice documents, developed in cooperation with Kuehne+Nagel, one of the world’s largest logistics companies. The primary task is to extract transport reference number from these documents. The dataset contains 764 valid invoice document and transport reference pairs. While the real world dataset contains sensitive customer data and cannot be released publicly, we provide an anonymized subset of the dataset and include the evaluation results in appendix.

The dataset presents several challenging characteristics that make it an ideal testbed for evaluating iterative learning capabilities. First, it requires specialized domain knowledge of business documents and terminology not commonly found in general language model training. Second, the hierarchical structure of UBL documents and the significant variability in format and identification patterns pose substantial extraction challenges. Additionally, as a novel benchmark without prior literature coverage, this dataset offers unique opportunities to assess agents’ adaptive learning abilities in a practical, high-stakes business context.

Evaluation Protocol. We report the success rate of the agent. The agent system is required to output the final result into a standardized format. The task is considered success when the output transport reference exactly matches the ground truth label. Given the diverse range of possible edge cases in the outputs, manually creating an extraction and comparison module may not cover all scenarios. Therefore, we employ a LLM judge (Zheng et al. 2023a) to compare the agent’s output with the ground truth.

Baselines for Comparison. We compare performance with previous works, including prompting LLM using Chain of Thought (Wei et al. 2022), Reflection (Shinn et al. 2024), and agent system with no memory (Wu et al. 2023). For a detailed description of the implementation details of the baselines, please refer to appendix C.

4.2 Main Results

We randomly select 60 samples from the dataset for training, the remaining 704 samples are reserved for testing performance. The maximum number of conversational turns between the assistant agent and the user proxy was capped at 5. The batch size was set to 14. That is, for each epoch, 14 tasks are sampled from training set and used for optimization. Since individual trajectories can become excessively lengthy, making the input to optimizer LLM larger than its

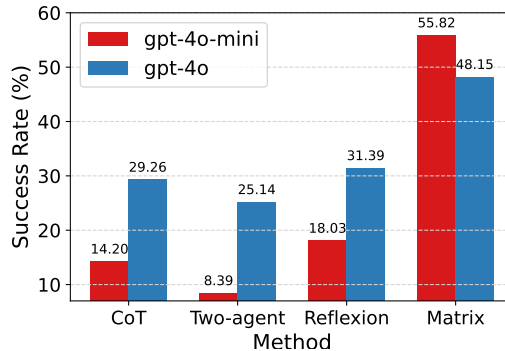


Figure 2: Comparison between Matrix and baselines with gpt-4o and gpt-4o-mini as backbone model. Matrix leverages gpt-4o for optimization in both cases. Surprisingly, gpt-4o-mini performs better after optimization.

context limit, we addressed this by truncating the batch to include only the largest number of trajectories that fit within the context limit. We set the backbone of the optimizer to gpt-4o. To maintain a fair comparison, we use gpt-4o for providing verbal feedback while running reflexion.

As shown in Figure 2, Matrix performs the best across all previous benchmarks. Without the memory module, vanilla agent cannot perform well without any prior knowledge, even performing worse than directly prompting the language model with chain-of-thought. In this case, equipping the agent with code interpreter actually hampers the performance, as the agent will over rely on writing code to extract strings rather than using natural language to reason. With the proposed Matrix mechanism, the agent’s task solving attempt is correctly guided by the long-term memory. The performance nearly doubles and outperforms all other methods.

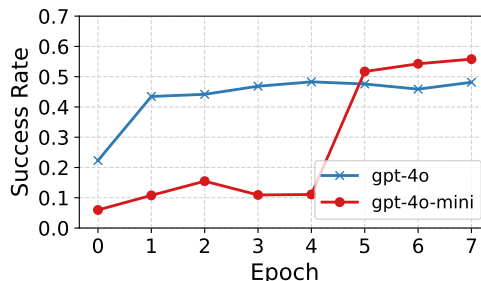


Figure 3: Success rate comparison between agent with gpt-4o and gpt-4o-mini as backbone over epochs.

Figure 3 demonstrates the performance of the agent with respect to the number of optimization epochs. The performance shows a steady rise across all epochs. With no memory, the agent system’s performance relies entirely on the capability of backbone LLM.

²<https://github.com/ag2ai/ag2>

4.3 Analysis of the Optimized System

In this section, we analyze the metrics of the system across the optimization of memory across epochs. Specifically, we explore three metrics: (1) average api calls it takes to correctly solve a question. This affects the overall latency of the agent system. (2) average cost it takes to correctly solve a question. (3) the distribution of the length of the successfully solved documents.

Analysis of the average number of API calls. As shown in Figure 4, the agent system exhibits a notable decrease in the average number of API calls required to solve the document question-answering task. After equipping with the optimized memory, the average number of API calls reduces about 8.12% for `gpt-4o` backed agent, 21.3% for `gpt-4o-mini` backed agent. This reduction indicates that the system is able to handle the task more efficiently with fewer calls, reducing sources of latency but reaching a better performance.

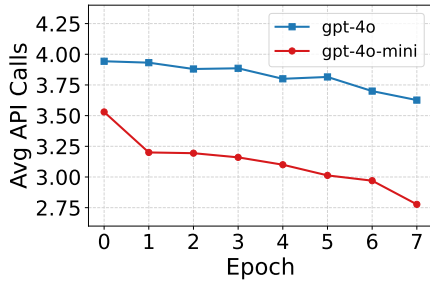


Figure 4: Comparison of average number of API calls it takes to solve a task. The average number decreases steadily as the training goes on.

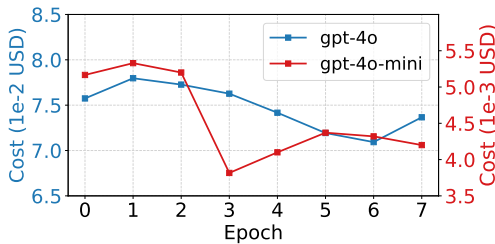


Figure 5: Average cost of API calls after each epoch. The cost shows a decreasing trend as training goes on.

Analysis of the cost. We plot the average cost of successfully solving a task after each epoch in Figure 5. The cost increases during the early stages of optimization. This rise is attributed to the initial memory initialization and optimization phase. During this phase, while the memory is introduced to guide problem-solving, its content is not yet fully refined. As a result, the agent writes more lengthy code for trial and subsequently leading to higher API costs. As optimization goes on, the memory gets gradually refined and more comprehensive, the agent can solve the task more efficiently with less tokens. The slight uptick in at the end of

the optimization suggests an adaptive process, where some rebalancing occurs after significant cost reduction. From a broader perspective, the overall trajectory suggests that optimization effectively lowers the long-term computational cost in solving tasks.

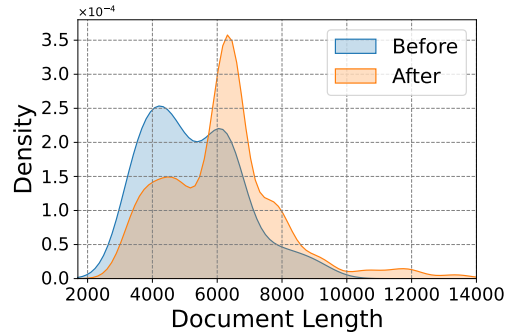


Figure 6: Distribution of successfully analyzed document lengths before and after optimization with `gpt-4o` as backbone.

Analysis of the length of successfully solved documents. As shown in Figure 6, the distribution of the successfully analyzed document lengths by `gpt-4o` backed agent is illustrated for two different stages of the agent’s performance: before optimization and after optimization. The document length is measured by the total number of tokens after tokenization. The distribution of successfully analyzed document lengths shifts notably after optimization. Before training, the performance peaks for document lengths around 4000 and 6000 tokens, and significantly declines after 6000 tokens. After training, we observe a significant improvement, especially for longer documents, with the distribution peak shifting towards 6,000 tokens and an extended tail that stretches beyond. This pattern showcases the enhanced capacity of the agent system to handle larger documents after optimization, reflecting its improved robustness and processing efficiency for more complex and lengthy inputs. We observe similar pattern for `gpt-4o-mini` backed agent and therefore omit the plotting for brevity.

4.4 Optimization with Weaker Language Model

In the main experiment, we leveraged `gpt-4o` as backbone for optimization. We next explore the performance with smaller language model used as optimizer. Following the protocol in Section 4.2, we set the backbone of optimizer and agents both to `gpt-4o-mini`.

Figure 7 demonstrates how the performance changes with respect to the number of training epochs. Although `gpt-4o-mini` has weaker language understanding capability, it can still understand the task and summarize reusable patterns from the task trajectories. However, it misses some regular expression patterns, leading to a less comprehensive optimized memory. Therefore, the final performance of the optimized system remains limited.

We then compare the performance between Matrix and Reflexion with `gpt-4o-mini` for optimization and pro-

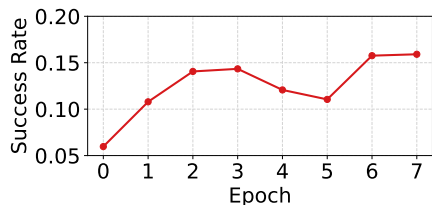


Figure 7: Performance of the optimized agent across iterations using `gpt-4o-mini` as optimizer. The performance rises in the beginning, showing steady improvement in success rate over the first three epochs, but soon begins to stabilize and reaches a slightly higher plateau.

Opt. LM	Reflexion	Matrix
<code>gpt-4o-mini</code>	12.64	15.90
<code>gpt-4o</code>	18.03	55.82

Table 1: Comparison between Reflexion and Matrix with `gpt-4o-mini` and `gpt-4o` as backbone for verbal feedback and optimization. In both cases, Matrix performs better than Reflexion.

viding verbal feedback, since they are both learning-based methods. The results, presented in Table 1, indicate that Matrix consistently outperforms Reflexion. Notably, the online algorithmic nature of Reflexion necessitates collecting verbal feedback for agents through trial and error on a task-by-task basis, which is inefficient as the number of test tasks increases. In contrast, Matrix leverages a generalized pattern learned from trajectories, enabling it to enhance performance in a more cost-efficient manner.

5 Conclusion

This paper explores the application of LLM agents for business document information retrieval, focusing on extracting transport references from invoices. To specialize LLM agents for this domain, we introduce Matrix, a paradigm enabling systematic learning through iterative self-refinement and memory updates. Real-world evaluations show Matrix distills actionable insights, outperforms baselines by large margin, and improves latency, cost, and document processing capability. This work reveals the potential of LLM agents for efficient, scalable enterprise document automation.

References

Abujabal, A.; Saha Roy, R.; Yahya, M.; and Weikum, G. 2019. ComQA: A Community-sourced Dataset for Complex Factoid Question Answering with Paraphrase Clusters. In Burstein, J.; Doran, C.; and Solorio, T., eds., *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 307–317. Minneapolis, Minnesota: Association for Computational Linguistics.

Achiam, J.; Adler, S.; Agarwal, S.; Ahmad, L.; Akkaya, I.; Aleman, F. L.; Almeida, D.; Altschmidt, J.; Altman, S.; Anadkat, S.; et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Antonacopoulos, A.; Bridson, D.; Papadopoulos, C.; and Pletschacher, S. 2009. A realistic dataset for performance evaluation of document layout analysis. In *2009 10th International Conference on Document Analysis and Recognition*, 296–300. IEEE.

Bajaj, P.; Campos, D.; Craswell, N.; Deng, L.; Gao, J.; Liu, X.; Majumder, R.; McNamara, A.; Mitra, B.; Nguyen, T.; et al. 2016. Ms marco: A human generated machine reading comprehension dataset. *arXiv preprint arXiv:1611.09268*.

Brown, T. B. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.

Cai, T.; Wang, X.; Ma, T.; Chen, X.; and Zhou, D. 2023. Large language models as tool makers. *arXiv preprint arXiv:2305.17126*.

Cheng, J.; Liu, X.; Zheng, K.; Ke, P.; Wang, H.; Dong, Y.; Tang, J.; and Huang, M. 2023. Black-box prompt optimization: Aligning large language models without model training. *arXiv preprint arXiv:2311.04155*.

Dong, Q.; Li, L.; Dai, D.; Zheng, C.; Wu, Z.; Chang, B.; Sun, X.; Xu, J.; and Sui, Z. 2022. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*.

Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; Al-Dahle, A.; Letman, A.; Mathur, A.; Schelten, A.; Yang, A.; Fan, A.; et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Guo, X.; Huang, K.; Liu, J.; Fan, W.; Vélez, N.; Wu, Q.; Wang, H.; Griffiths, T. L.; and Wang, M. 2024. Embodied llm agents learn to cooperate in organized teams. *arXiv preprint arXiv:2403.12482*.

Hamdi, A.; Carel, E.; Joseph, A.; Coustaty, M.; and Doucet, A. 2021. Information extraction from invoices. In *International Conference on Document Analysis and Recognition*, 699–714. Springer.

Harley, A. W.; Ufkes, A.; and Derpanis, K. G. 2015. Evaluation of deep convolutional nets for document image classification and retrieval. In *2015 13th International Conference on Document Analysis and Recognition (ICDAR)*, 991–995. IEEE.

Hu, S.; Lu, C.; and Clune, J. 2024. Automated design of agentic systems. *arXiv preprint arXiv:2408.08435*.

Huang, J.; Chen, X.; Mishra, S.; Zheng, H. S.; Yu, A. W.; Song, X.; and Zhou, D. 2023. Large language models cannot self-correct reasoning yet. *arXiv preprint arXiv:2310.01798*.

Ji, Z.; Yu, T.; Xu, Y.; Lee, N.; Ishii, E.; and Fung, P. 2023. Towards mitigating LLM hallucination via self reflection. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, 1827–1843.

Jiang, D.; Zhang, J.; Weller, O.; Weir, N.; Van Durme, B.; and Khashabi, D. 2024. Self-[in] correct: LLMs struggle with refining self-generated responses. *arXiv preprint arXiv:2404.04298*.

- Joshi, M.; Choi, E.; Weld, D.; and Zettlemoyer, L. 2017. TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension. In Barzilay, R.; and Kan, M.-Y., eds., *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1601–1611. Vancouver, Canada: Association for Computational Linguistics.
- Kamoi, R.; Zhang, Y.; Zhang, N.; Han, J.; and Zhang, R. 2024. When Can LLMs Actually Correct Their Own Mistakes? A Critical Survey of Self-Correction of LLMs. *arXiv preprint arXiv:2406.01297*.
- Krieger, F.; Drews, P.; Funk, B.; and Wobbe, T. 2021. Information extraction from invoices: a graph neural network approach for datasets with high layout variety. In *Innovation Through Information Systems: Volume II: A Collection of Latest Research on Technology Issues*, 5–20. Springer.
- Kwiatkowski, T.; Palomaki, J.; Redfield, O.; Collins, M.; Parikh, A.; Alberti, C.; Epstein, D.; Polosukhin, I.; Devlin, J.; Lee, K.; et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7: 453–466.
- Liu, J.; Shen, D.; Zhang, Y.; Dolan, B.; Carin, L.; and Chen, W. 2021. What Makes Good In-Context Examples for GPT-3? *arXiv preprint arXiv:2101.06804*.
- Lu, Y.; Bartolo, M.; Moore, A.; Riedel, S.; and Stenetorp, P. 2021. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv preprint arXiv:2104.08786*.
- Madaan, A.; Tandon, N.; Gupta, P.; Hallinan, S.; Gao, L.; Wiegrefe, S.; Alon, U.; Dziri, N.; Prabhunoye, S.; Yang, Y.; et al. 2024. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36.
- Min, S.; Lewis, M.; Zettlemoyer, L.; and Hajishirzi, H. 2021. Metaicl: Learning to learn in context. *arXiv preprint arXiv:2110.15943*.
- Palm, R. B.; Winther, O.; and Laws, F. 2017. Cloudscan—a configuration-free invoice analysis system using recurrent neural networks. In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, volume 1, 406–413. IEEE.
- Pan, L.; Saxon, M.; Xu, W.; Nathani, D.; Wang, X.; and Wang, W. Y. 2023. Automatically correcting large language models: Surveying the landscape of diverse self-correction strategies. *arXiv preprint arXiv:2308.03188*.
- Prasad, A.; Hase, P.; Zhou, X.; and Bansal, M. 2022. Grips: Gradient-free, edit-based instruction search for prompting large language models. *arXiv preprint arXiv:2203.07281*.
- Pryzant, R.; Iter, D.; Li, J.; Lee, Y. T.; Zhu, C.; and Zeng, M. 2023. Automatic prompt optimization with “gradient descent” and beam search. *arXiv preprint arXiv:2305.03495*.
- Qian, C.; Han, C.; Fung, Y. R.; Qin, Y.; Liu, Z.; and Ji, H. 2023. Creator: Tool creation for disentangling abstract and concrete reasoning of large language models. *arXiv preprint arXiv:2305.14318*.
- Rajpurkar, P.; Zhang, J.; Lopyrev, K.; and Liang, P. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. In Su, J.; Duh, K.; and Carreras, X., eds., *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2383–2392. Austin, Texas: Association for Computational Linguistics.
- Riba, P.; Dutta, A.; Goldmann, L.; Fornés, A.; Ramos, O.; and Lladós, J. 2019. Table detection in invoice documents by graph neural networks. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, 122–127. IEEE.
- Shinn, N.; Cassano, F.; Gopinath, A.; Narasimhan, K.; and Yao, S. 2024. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36.
- Song, L.; Liu, J.; Zhang, J.; Zhang, S.; Luo, A.; Wang, S.; Wu, Q.; and Wang, C. 2024. Adaptive In-conversation Team Building for Language Model Agents. *arXiv preprint arXiv:2405.19425*.
- Su, H.; Kasai, J.; Wu, C. H.; Shi, W.; Wang, T.; Xin, J.; Zhang, R.; Ostendorf, M.; Zettlemoyer, L.; Smith, N. A.; et al. 2022. Selective annotation makes language models better few-shot learners. *arXiv preprint arXiv:2209.01975*.
- Sumers, T. R.; Yao, S.; Narasimhan, K.; and Griffiths, T. L. 2023. Cognitive architectures for language agents. *arXiv preprint arXiv:2309.02427*.
- Talmor, A.; and Berant, J. 2018. The Web as a Knowledge-base for Answering Complex Questions. In *North American Association for Computational Linguistics (NAACL)*.
- Tarawneh, A. S.; Hassanat, A. B.; Chetverikov, D.; Lendak, I.; and Verma, C. 2019. Invoice classification using deep features and machine learning techniques. In *2019 IEEE Jordan International Joint Conference on Electrical Engineering and Information Technology (JEEIT)*, 855–859. IEEE.
- Wang, C.; Liu, X.; Yue, Y.; Tang, X.; Zhang, T.; Jiayang, C.; Yao, Y.; Gao, W.; Hu, X.; Qi, Z.; et al. 2023. Survey on factuality in large language models: Knowledge, retrieval and domain-specificity. *arXiv preprint arXiv:2310.07521*.
- Wang, C.; Wu, Q.; Weimer, M.; and Zhu, E. 2021. Flam: A fast and lightweight automl library. *Proceedings of Machine Learning and Systems*, 3: 434–447.
- Wang, Z. Z.; Mao, J.; Fried, D.; and Neubig, G. 2024. Agent workflow memory. *arXiv preprint arXiv:2409.07429*.
- Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Xia, F.; Chi, E.; Le, Q. V.; Zhou, D.; et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35: 24824–24837.
- Wu, Q.; Bansal, G.; Zhang, J.; Wu, Y.; Zhang, S.; Zhu, E.; Li, B.; Jiang, L.; Zhang, X.; and Wang, C. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*.
- Wu, Y.; Jia, F.; Zhang, S.; Li, H.; Zhu, E.; Wang, Y.; Lee, Y. T.; Peng, R.; Wu, Q.; and Wang, C. 2024a. MathChat: Converse to Tackle Challenging Math Problems with LLM Agents. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*.

- Wu, Y.; Yue, T.; Zhang, S.; Wang, C.; and Wu, Q. 2024b. StateFlow: Enhancing LLM Task-Solving through State-Driven Workflows. *arXiv preprint arXiv:2403.11322*.
- Wu, Z.; Wang, Y.; Ye, J.; and Kong, L. 2022. Self-adaptive in-context learning: An information compression perspective for in-context example selection and ordering. *arXiv preprint arXiv:2212.10375*.
- Xia, X.; Liu, J.; Yu, J.; Shen, X.; Han, B.; and Liu, T. 2022. Moderate coreset: A universal method of data selection for real-world data-efficient deep learning. In *The Eleventh International Conference on Learning Representations*.
- Xia, X.; Liu, J.; Zhang, S.; Wu, Q.; Wei, H.; and Liu, T. 2024. Refined Coreset Selection: Towards Minimal Coreset Size under Model Performance Constraints. In *Forty-first International Conference on Machine Learning*.
- Yang, A.; Yang, B.; Hui, B.; Zheng, B.; Yu, B.; Zhou, C.; Li, C.; Li, C.; Liu, D.; Huang, F.; et al. 2024. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*.
- Yang, C.; Wang, X.; Lu, Y.; Liu, H.; Le, Q. V.; Zhou, D.; and Chen, X. 2023. Large language models as optimizers. *arXiv preprint arXiv:2309.03409*.
- Yang, Z.; Qi, P.; Zhang, S.; Bengio, Y.; Cohen, W. W.; Salakhutdinov, R.; and Manning, C. D. 2018. HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Yao, S.; Yu, D.; Zhao, J.; Shafran, I.; Griffiths, T.; Cao, Y.; and Narasimhan, K. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.
- Yao, S.; Zhao, J.; Yu, D.; Du, N.; Shafran, I.; Narasimhan, K.; and Cao, Y. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. In *International Conference on Learning Representations (ICLR)*.
- Yuan, L.; Chen, Y.; Wang, X.; Fung, Y. R.; Peng, H.; and Ji, H. 2023. Craft: Customizing llms by creating and retrieving from specialized toolsets. *arXiv preprint arXiv:2309.17428*.
- Zhang, J.; Xiang, J.; Yu, Z.; Teng, F.; Chen, X.; Chen, J.; Zhuge, M.; Cheng, X.; Hong, S.; Wang, J.; et al. 2024a. AFlow: Automating Agentic Workflow Generation. *arXiv preprint arXiv:2410.10762*.
- Zhang, S.; Jia, F.; Wang, C.; and Wu, Q. 2023a. Targeted hyperparameter optimization with lexicographic preferences over multiple objectives. In *The Eleventh international conference on learning representations*.
- Zhang, S.; Wu, Y.; Zheng, Z.; Wu, Q.; and Wang, C. 2024b. Hypertime: Hyperparameter optimization for combating temporal distribution shifts. In *Proceedings of the 32nd ACM International Conference on Multimedia*, 4610–4619.
- Zhang, S.; Xia, X.; Wang, Z.; Chen, L.-H.; Liu, J.; Wu, Q.; and Liu, T. 2023b. Ideal: Influence-driven selective annotations empower in-context learners in large language models. *arXiv preprint arXiv:2310.10873*.
- Zhang, S.; Zhang, J.; Liu, J.; Song, L.; Wang, C.; Krishna, R.; and Wu, Q. 2024c. Offline Training of Language Model Agents with Functions as Learnable Weights. In *Forty-first International Conference on Machine Learning*.
- Zhao, A.; Huang, D.; Xu, Q.; Lin, M.; Liu, Y.-J.; and Huang, G. 2024. Expel: Llm agents are experiential learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 19632–19642.
- Zhao, Z.; Wallace, E.; Feng, S.; Klein, D.; and Singh, S. 2021. Calibrate Before Use: Improving Few-shot Performance of Language Models. In Meila, M.; and Zhang, T., eds., *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, 12697–12706. PMLR.
- Zheng, L.; Chiang, W.-L.; Sheng, Y.; Zhuang, S.; Wu, Z.; Zhuang, Y.; Lin, Z.; Li, Z.; Li, D.; Xing, E.; et al. 2023a. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36: 46595–46623.
- Zheng, X.; Yang, C.; Zhang, S.; Wang, Y.; Zhang, B.; Wu, Y.; Wu, Y.; Shao, L.; and Ji, R. 2023b. Ddpnas: Efficient neural architecture search via dynamic distribution pruning. *International Journal of Computer Vision*, 131(5): 1234–1249.
- Zhong, X.; Tang, J.; and Yepes, A. J. 2019. Publaynet: largest dataset ever for document layout analysis. In *2019 International conference on document analysis and recognition (ICDAR)*, 1015–1022. IEEE.
- Zhou, D.; Schärli, N.; Hou, L.; Wei, J.; Scales, N.; Wang, X.; Schuurmans, D.; Cui, C.; Bousquet, O.; Le, Q.; et al. 2022a. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*.
- Zhou, Y.; Muresanu, A. I.; Han, Z.; Paster, K.; Pitis, S.; Chan, H.; and Ba, J. 2022b. Large language models are human-level prompt engineers. *arXiv preprint arXiv:2211.01910*.
- Zhuge, M.; Wang, W.; Kirsch, L.; Faccio, F.; Khizbullin, D.; and Schmidhuber, J. 2024. Language agents as optimizable graphs. *arXiv preprint arXiv:2402.16823*.

A Dataset Details

The dataset we experiment on contains sensitive data that reveals the customers’ business details. Therefore, we are unable to release the original dataset. We will release an anonymized subset of the dataset to facilitate future research in this field. The goal of anonymization is to preserve the original data structure while removing sensitive information. In the following section, we describe the anonymization pipeline.

A.1 Background and data source

The dataset used in this study originates from real-world invoices processed by Kuehne+Nagel. These invoices represent a diverse range of business transactions, providing a rich source of data for our transport reference and tracking number processing task.

A.2 Data collection process

The data collection was facilitated through the Beyond Work platform, a human/AI collaboration platform for solving tedious work. Specifically, the invoices were processed within a workblock, which is a container for all technology (code, infrastructure, permission etc) needed to solve a distinct process.

The workblock used to collect this data, targets the process of ”Parked Invoices” in which invoices that cannot be automatically processed have to be processed by humans. The workblock uses LLMs to process the invoice and introduces a validation task in which human workers were tasked with validating tracking numbers and transport references extracted from the invoices. This process not only involved identification and verification of the correct information but also required workers to provide explanations for any discrepancies or errors they encountered. This human-in-the-loop approach ensured high-quality, validated data.

A.3 Data preparation and anonymization

Following the validation process, the data underwent several stages of preparation to ensure its suitability for research purposes while maintaining strict privacy standards. The anonymization of the data included both pseudonymization of the identifiers we were looking to extract for the invoices, as the original human validation of the identifiers was done on the non-anonymized data. This way we ensure to keep the format of the identifier while removing any traceability to the original numbers. For all other sensitive data we completely anonymized it. A separate workblock was authored (the process of using natural language to create workblocks) specifically for the anonymization process. This step was crucial to protect sensitive business information and comply with data privacy regulations.

A.4 Unstructured data pseudonymization

For unstructured text data, we specifically used the claude-3.5-Sonnet language model within the anonymization workblock. This helped us in replacing identifying information with pseudonyms while maintaining the contextual and structural integrity of the data.

A.5 Structured data anonymization

For more structured data fields, we applied complete anonymization. This process involved replacing text data with random strings, integer values with random integers, and float values with random floats. This approach ensures that no traceable information remains in the structured fields.

A.6 Resulting dataset

The resulting dataset retains the complex structure and challenges of real-world invoice data while being fully anonymized. It preserves the intricate nature of business documents, allowing for realistic evaluation of information extraction techniques while ensuring the confidentiality of the original data sources.

B Results on Anonymized Data

In this section, we present the main results of the proposed method on the released anonymized data. The dataset comprises 127 task instances, with each instance consisting of an invoice document paired with its corresponding transport reference. In real-world scenarios, an invoice may not always contain a valid transport reference, and this characteristic has been intentionally preserved in the anonymized dataset to reflect real-world conditions. Of the 127 documents, 50 contain valid transport references, which form the primary focus of our study³.

Our work specifically focuses on cases where transport references are present. While we acknowledge the limited size of the anonymized dataset, this constraint arises from the significant manual effort required for anonymization and rigorous inspection. These measures are essential to ensure compliance with data privacy standards and to prevent the leakage of sensitive

³Data available at <https://tinyurl.com/3pakk9t4>

information. Our future work involves expanding the anonymized dataset to accommodate scenarios with missing or invalid transport references to provide a broader evaluation of our method.

We randomly select 8 task instances for training, 42 for testing performance. Due to the training set is much smaller than that of the main experiment, the agent’s task solving trajectory and the reflections can fit into the optimizer’s context window. In the training phase, the agent repeatedly attempts the same batch of tasks and the memory gets updated accordingly.

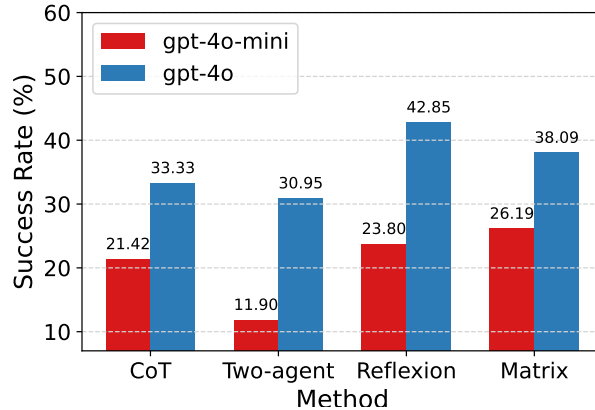


Figure 8: Performance comparison of baselines and Matrix on anonymized dataset. Reflexion performs better than Matrix on gpt-4o agent. In the remaining cases, Matrix performs better.

We present the performance results on the anonymized dataset in Figure 8. While Reflexion outperforms Matrix when using gpt-4o as the backbone, Matrix achieves the best performance across all other configurations. However, we observe that Matrix does not exhibit significant improvements over other baseline methods on the anonymized dataset. This is due to the limitations of the training dataset, which is not sufficiently large or diverse to provide a well-generalized representation of the problem space. Specifically, the optimizer is constrained to use the same batch of training tasks across all epochs, resulting in a lack of diversity in the gathered task solving experience. Consequently, the optimizer struggles to capture a comprehensive pattern that accurately reflects the full data distribution.

This highlights a key limitation of Matrix: it requires a substantial amount of training data to effectively model and generalize the patterns within the dataset. Addressing the challenge of agent training under limited data resource constraints remains an avenue for future research.

C Implementation Details of baselines

C.1 Chain-of-thought

CoT is often used as the default way of prompting LLM. We follow this guideline and prompt the LLM to think step by step and analyze the dataset before generating the final answer.

C.2 Two-agent

The vanilla two agent system is implemented by AG2⁴ (formerly AutoGen) (Wu et al. 2023). In the system, an assistant agent backed by LLM is responsible for analyzing the task environment, reason, and make decisions. A user proxy agent receives the content generated by assistant agent, execute the python code provided, and automatically provides the feedback. For a detailed description of the system, please refer to the official documentation⁵.

C.3 Reflexion

Reflexion proposes to improve language agents using verbal feedback. For each task, the agent reflects on their performance and store these reflections in memory to make better decisions in the next trial. The system maintains specific memories for each individual task or instance. We change the ReAct agent of the original paper into the two-agent system we investigate, and run reflection on each task. The maximum number of trials is set to 7 to ensure a fair comparison. The prompt for generating reflection is identical to that of their original Github repository⁶.

⁴<https://github.com/ag2ai/ag2>

⁵<https://ag2ai.github.io/ag2/docs/tutorial/introduction>

⁶https://github.com/noahshinn/reflexion/blob/main/alfworld_runs/generate_reflections.py#L15

D Limitations and future work

As shown in Section 4.4, the performance of Matrix strongly depends on the capability of backbone models. Weaker models cannot distill actionable insights from experiences and correctly guide future task solving attempts.

The process of agent training requires a number of task instances that are representative of the full data distribution. Comparing performance of Matrix on full data and the anonymized data, we discover that the method requires larger training data to reach strong performance. The anonymized data contains little samples that cannot fully reflect the dataset distribution, therefore Matrix tends to underperform. How to identify a subset of the most important and influential samples (i.e. coresets selection (Xia et al. 2022, 2024)) for training is an open question for agent training. One potential solution is to leverage AutoML based methods (Wang et al. 2021; Zheng et al. 2023b; Zhang et al. 2023a, 2024b).