# Positive Experience Reflection for Agents in Interactive Text Environments

Philip Lippmann Matthijs T.J. Spaan Jie Yang

Delft University of Technology {p.lippmann, m.t.j.spaan, j.yang-3}@tudelft.nl

#### Abstract

Intelligent agents designed for interactive environments face significant challenges in text-based games, a domain that demands complex reasoning and adaptability. While agents based on large language models (LLMs) using self-reflection have shown promise, they struggle when initially successful and exhibit reduced effectiveness when using smaller LLMs. We introduce Sweet&Sour, a novel approach that addresses these limitations in existing reflection methods by incorporating positive experiences and managed memory to enrich the context available to the agent at decision time. Our comprehensive analysis spans both closed- and opensource LLMs and demonstrates the effectiveness of Sweet&Sour in improving agent performance, particularly in scenarios where previous approaches fall short.

## **1** Introduction

Intelligent agents, designed to interact with and make decisions in dynamic environments, have become a central focus in AI research, with text-based games (TBGs) emerging as a particularly challenging domain for evaluating these agents' reasoning, adaptability, and learning abilities [2, 21]. Originally popular in the 1970s as text adventure games,<sup>1</sup> TBGs present players with textual descriptions of environments, requiring them to input natural language commands to achieve objectives [4]. For instance, determining if a metal fork is conductive involves locating the fork, assembling a circuit, and analyzing the result. Navigating TBGs demands that agents exhibit a combination of abilities, including planning, memory retention, spatial reasoning, and common sense knowledge [22].

Previously, deep reinforcement learning and behavior cloning were the primary approaches to develop agents to play TBGs [1, 25]. However, recent research shows that agents based on pretrained large language models (LLMs) are more effective at navigating TBGs [11]. A key factor in their success is the integration of internal *reflection* to improve planning [24, 9, 7].

Self-reflection, closely related to self-refinement, is a form of reasoning that occurs after receiving binary or scalar feedback from the environment [12]. In this process, the LLM reviews its actions and their outcomes, considering what went wrong and potential ways to improve [23]. By iteratively adjusting its strategy based on verbal reinforcement, conveyed through textual feedback, the agent refines its planning for subsequent attempts [17]. However, reflection also has several limitations, including 1) underwhelming performance when agents are correct initially [10], 2) significantly worse efficacy when using smaller LLMs [11], and 3) dependence on external feedback [29].

**Our Contributions** In this work we conduct a comprehensive analysis of LLM-based agents employing reflection approaches in TBGs and evaluate their performance across closed- and open-source LLMs. To address the limitations of poor performance when agents are initially successful and the diminished efficacy of smaller LLMs, we propose *Sweet&Sour* to leverage positive experiences and

<sup>&</sup>lt;sup>1</sup>Try it yourself: https://www.microsoft.com/en-us/research/project/textworld/try-it/

<sup>38</sup>th Conference on Neural Information Processing Systems (NeurIPS 2024).



Figure 1: Comparison of used prompting methods to play ScienceWorld. ReAct introduces a THINK action to explicitly reason regarding the next step. Reflexion leverages self-reflection across attempts to learn from unsuccessful tries and stores these in memory. Sweet&Sour not only performs self-reflection after failures but also after each completed sub goal, making its reflection instantly available.

managed memory to create a richer context for self-reflection. Our findings demonstrate that our method improves the performance of agents using reflection, particularly in scenarios where they previously struggled, enabling more robust and generalizable learning across tasks and model sizes.

# 2 Methodology

**Background** Assuming an LLM behaving as an actor model as part of our agent – *i.e.* generating actions based on the current state and policy, analogous to traditional policy-based RL setups – we sample an action  $a_t$  from the current policy  $\pi_{\theta}$  at time t and receive an observation from the environment  $o_t$ . When a game begins, the agent makes its first observation  $o_0$  at time step t = 0. This first observation differs from subsequent ones, as it consists of the goal description d, as well as an analysis of the starting room (*i.e.* the output of the "look around" command). Subsequently, the agent can perform an action  $a_t$  at each time step and receives a corresponding observation  $o_t$  from the environment. The agent has an inventory  $i_t$  in which to store items. Each task consists of a number of sub tasks (such as finding a key object), the completion of which grants the agent a sparse reward, which adds to its current reward  $r_t$ . The game continues until the agent has achieved the goal outlined in d and receives the full reward as final score, or the maximum number of steps – which we set to 150 – is reached, in which case  $r_t$  will become the final score. A detailed problem formulation and assumptions we make are given in appendix A.

**Self-Reflection** Reflection occurs in addition to the acting LLM. Here, the agent reviews the  $a_t$  and  $o_t$  associated with previous unsuccessful attempts to verbalize the reason for failure. This process typically involves maintaining a persistent history of insights gained across attempts, which the LLM uses as additional context for its reflections to improve future decision making for the next attempt [17]. However, since other self-reflection methods focus on learning from failures [16, 28, 30, 8, 27], they overlook the importance of reinforcing successful behaviors in a similar way.

**Sweet&Sour** To address the limitations of existing reflection methods, we enhance self-reflection by incorporating both positive (*sweet*) and negative (*sour*) experiences into the reflection process. This enables the agent to learn effectively from successful actions by reinforcing strategies that lead to positive outcomes while still learning from failures. When the current policy is achieving rewards, we query the agent to extrapolate from it, encouraging the agent to verbalize what made its current policy successful and what can be generalized from this. This is visualized in figure 1 and an example of this is shown in appendix B. Our method is broadly applicable to agents in interactive text environments

with feedback using self-reflection, including those that build additional complexity on top of the core reflection loop, such as grounding [11] or gradient learning [27]. We leave the study of these additional use cases for future work.

Previous works store their reflections gained from unsuccessful attempts in something akin to a longterm memory to make them available to the agent across attempts [8, 17]. Instead, to complement Sweet&Sour, we propose a *managed memory* approach to store and retrieve relevant reflections. This is implemented using a dual-buffer structure, where experiences are stored in two categories: short-term memory and long-term memory, based on their outcome (success or failure) and recency. Initially, if sub goals are reached, short-term memories of the made reflections are stored in a temporary buffer. Each short-term memory consists of a tuple (reflection<sub>t</sub>,  $o_t$ ,  $a_t$ ,  $r_t$ ). Once a task is completed or an attempt ends, all short-term memories are moved to long-term memory. Failed attempts' reflections are immediately added to long-term memory for the next attempt and short-term memory collection ends.

# **3** Experiments

**Data & Environment** We use the ScienceWorld benchmark [21], which provides a versatile setting for evaluating agents in science experiment tasks across 10 interconnected locations, such as a greenhouse and a workshop, with over 200 objects and 25 action templates, generating a vast and dynamic search space. We use the test set for our evaluation, which provides up to 10 variations of each of the 30 distinct tasks. These interactive tasks cover various topics, including chemistry and electricity, and have an average optimal decision depth of 50 steps. An example task is shown in appendix C. For details of all tasks and the environment, we refer to [21]. We elect to use ScienceWorld instead of previous interactive text environment benchmarks such as TWC [13] and ALFWorld [18], due to their relative simplicity for current LLM-based agents. We measure performance using the success score, which is always between 0 and 100. Completing a task implies completing every sub task, receiving the full reward, and thus a score of 100.

**Baselines** CALM [25] is a reranking method that integrates a deep reinforced relevance network (DRRN) [5] with a causal language model fine-tuned using oracle trajectories. The causal language model acquires task- and environment-specific knowledge through imitation learning, while DRRN trains a policy network to rerank the language model's predictions. We use ReAct [26] as our baseline LLM-based agent. ReAct composes useful information at each time step by reasoning over the current context (*e.g.* decomposing task or common sense knowledge query) and carries it forward to the context of the following time step. This reasoning action does not affect the environment and may be considered few-shot in-context learning. To contextualize our work, we compare our method against Reflexion [17], an agent built on ReAct that employs a self-reflection mechanism to iteratively improve its performance across rounds upon encountering failure based on feedback from the environment. As such, it runs over up to four rounds as it builds up its long-term memory. For all agents, we evaluate their performance using LLMs of different sizes and complexities to assess the performance of each method across varying computational resources. In descending order of parameter count, we select GPT-40 (gpt-40-2024-08-06)[15], Mistral Large 2 (mistral-large-2407), and Llama 3.1 8B (llama-3.1-8b-instruct)[3], accessing each through its respective APIs.

**Main Results** The results are shown in table 1. We find that Sweet&Sour outperforms the baseline methods across all LLMs, setting the highest average score at 54.6 using GPT-40. The performance gap between Sweet&Sour and the other methods widens for smaller models with a lower parameter count. For instance, it achieves 44.6 compared to Reflexion's 27.6 on Mistral Large 2, and 32.5 compared to 21.7 on Llama 8B – indicating that our method is more suitable for scenarios with limited computational resources. When we modify our method to only sample from failures, performance drops significantly to a level similar to Reflexion – scores decrease to 24.6, 31.1, and 44.9 for Llama 8B, Mistral Large 2, and GPT-40, respectively. As such, incorporating positive experiences indeed leads to better reflections, mimicking how humans learn from both positive and negative experiences, resulting in improved decision-making.

**Anti-Tilt** In highly challenging tasks, such as 1-1 and 8-2, all methods tend to struggle, while in simpler tasks, most methods succeed based on the LLM's inherent capabilities alone. However, medium-difficulty tasks, such as 3-2 and 3-3, reveal a critical performance gap between our method and previous approaches. We note that this occurs because traditional methods fail to reflect on early

Task	CALM	ReAct			Reflexion			Sweet&Sour (ours)		
	CALM	L8B	ML2	GPT	L8B	ML2	GPT	L8B	ML2	GPT
1-1 (Boil)	0.0	0.0	0.0	3.8	0.0	0.0	5.1	0.0	7.2	9.6
1-2 (Melt)	0.0	8.4	10.3	11.8	0.0	0.0	10.0	11.4	12.1	12.8
1-3 (Freeze)	0.0	1.5	0.0	8.1	0.0	2.3	8.3	2.4	3.1	8.9
1-4 (Change state)	0.0	1.0	4.7	10.0	0.0	0.0	4.2	1.7	2.9	9.2
2-1 (Thermometer)	1.0	5.1	7.8	7.7	3.4	4.2	7.6	7.8	9.7	10.9
2-2 (Melting)	1.0	6.7	6.3	5.9	3.3	3.3	26.2	7.9	36.8	46.0
2-3 (Melting)	5.0	9.1	11.8	23.4	13.2	14.7	22.6	15.2	29.0	38.3
<b>3-1</b> (Power 1)	7.0	18.8	24.6	57.2	21.2	51.5	78.4	28.6	75.4	81.1
<b>3-2</b> (Power 2)	2.0	10.2	24.7	55.6	9.5	11.9	24.7	23.3	44.5	58.0
<b>3-3</b> (Conductivity 1)	2.0	52.4	51.7	73.0	9.2	25.8	72.1	59.1	69.2	75.7
<b>3-4</b> (Conductivity 2)	10.0	54.2	64.9	89.7	35.4	41.6	75.1	62.7	60.3	67.3
<b>4-1</b> (Find 1)	54.0	17.3	18.7	27.5	44.6	48.1	62.3	41.7	71.7	74.2
<b>4-2</b> (Find 2)	10.0	69.1	71.6	80.3	68.4	75.7	87.3	76.8	100.0	100.0
<b>4-3</b> (Find 3)	8.0	21.3	42.8	47.7	18.4	16.5	17.3	20.9	21.5	34.3
<b>4-4</b> (Find 4)	2.0	15.7	15.2	19.3	39.6	46.6	100.0	55.1	87.8	100.0
5-1 (Grow plant)	4.0	10.8	10.8	10.0	7.2	7.2	7.9	14.2	14.6	17.4
<b>5-2</b> (Grow fruit)	3.0	18.1	18.5	19.2	30.8	51.4	34.6	51.5	55.6	60.2
6-1 (Chemistry 1)	6.0	37.8	42.9	58.6	27.1	29.7	70.2	37.9	61.1	70.2
6-2 (Chemistry 2)	3.0	25.0	27.1	50.6	14.4	28.0	69.8	27.2	51.9	83.1
<b>6-3</b> (Chemistry 3)	6.0	14.4	17.5	39.7	38.9	31.1	16.7	45.3	53.7	61.5
7-1 (Lifespan 1)	10.0	37.0	41.7	60.0	75.0	75.0	100.0	75.0	88.2	100.0
7-2 (Lifespan 2)	4.0	50.5	50.7	67.5	60.0	71.9	81.4	70.5	77.0	80.0
<b>7-3</b> (Lifespan 3)	4.0	33.7	38.2	50.0	29.5	33.7	75.0	51.1	54.2	84.6
8-1 (Identify life 1)	0.0	5.1	18.9	25.3	1.7	1.7	3.4	11.1	10.3	14.2
8-2 (Identify life 2)	0.0	6.4	7.4	8.0	7.4	8.0	8.0	5.0	7.4	7.4
9-1 (Measure angle)	0.0	28.5	33.0	42.5	56.9	55.1	57.1	68.4	70.3	75.0
<b>9-2</b> (Friction 1)	3.0	14.5	22.6	43.1	23.4	29.3	100.0	33.3	36.7	62.0
<b>9-3</b> (Friction 2)	2.0	2.9	14.5	42.8	1.3	33.6	59.6	7.2	51.9	63.1
<b>10-1</b> (Genetics 1)	2.0	25.7	27.3	26.4	5.6	9.8	50.4	38.9	48.6	78.8
<b>10-2</b> (Genetics 2)	2.0	13.2	19.1	17.2	6.2	21.5	22.7	23.6	24.0	54.8
Average	5.07	20.5	24.8	36.0	21.7	27.6	45.3	32.5	44.6	54.6

Table 1: Results on the ScienceWorld benchmark. For each method, we use GPT-40 (GPT), Mistral Large 2 (ML2), and Llama 8B (L8B). Each value is an average of across all task variations.

successes, leaving them with less contextual understanding to carry momentum forward, leading to a sustained performance decline or "tilt." By contrast, we theorize that Sweet&Sour's reflection on both successes and failures provides a richer context, enabling it to sustain and build upon initial successes. For instance, on task 3-2, Sweet&Sour achieves a score of 68.0 with GPT-40, significantly outperforming Reflexion's 24.7 while beating ReAct's 55.6 by a smaller margin, indicating that Reflexion, although it is overall more capable, got stuck despite its self-reflection.

## 4 Conclusion

In an attempt to improve agent performance and reduce sensitivity to the LLM used, our study embarks on an analysis of reflection mechanisms for LLM-based agents. Existing work focuses on learning from failures only. Although these methods improve planning by analyzing past mistakes, they struggle when agents are initially successful and when using smaller LLMs. The primary contribution of our work is Sweet&Sour, a novel reflection technique for LLM-based agents in TBGs that leverages positive experiences to improve agent self-reflection. Our comprehensive analysis demonstrates the effectiveness of Sweet&Sour in enhancing agent adaptability and reasoning, particularly in challenging situations where previous approaches struggled.

**Limitations** Despite promising results, our work has limitations. LLMs do not provide guarantees regarding their reasoning capabilities. Additionally, our evaluation is conducted using a single environment, which, while comprehensive, does not cover all types of interactive scenarios. We leave the exploration of additional environments to future work.

## Acknowledgments

This work was supported by an Oracle for Research Grant Award, as well as SURF Grant EINF-8535.

#### References

- [1] Prithviraj Ammanabrolu and Mark Riedl. Playing text-adventure games with graph-based deep reinforcement learning. In 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), pages 3557–3565, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1358. URL https://aclanthology. org/N19-1358.
- [2] Marc-Alexandre Côté, Ákos Kádár, Xingdi Yuan, Ben Kybartas, Tavian Barnes, Emery Fine, James Moore, Ruo Yu Tao, Matthew Hausknecht, Layla El Asri, Mahmoud Adada, Wendy Tay, and Adam Trischler. Textworld: A learning environment for text-based games. *CoRR*, abs/1806.11532, 2018.
- [3] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet,

Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.

- [4] Matthew Hausknecht, Prithviraj Ammanabrolu, Marc-Alexandre Côté, and Xingdi Yuan. Interactive fiction games: A colossal adventure. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7903–7910, 2020.
- [5] Ji He, Jianshu Chen, Xiaodong He, Jianfeng Gao, Lihong Li, Li Deng, and Mari Ostendorf. Deep reinforcement learning with a natural language action space. In Katrin Erk and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1621–1630, Berlin, Germany, Au-

gust 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1153. URL https://aclanthology.org/P16-1153.

- [6] Ji He, Jianshu Chen, Xiaodong He, Jianfeng Gao, Lihong Li, Li Deng, and Mari Ostendorf. Deep reinforcement learning with a natural language action space. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1621–1630, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1153. URL https://aclanthology.org/P16-1153.
- [7] Sihao Hu, Tiansheng Huang, Fatih Ilhan, Selim Tekin, Gaowen Liu, Ramana Kompella, and Ling Liu. A survey on large language model-based game agents, 2024. URL https://arxiv. org/abs/2404.02039.
- [8] Xu Huang, Jianxun Lian, Yuxuan Lei, Jing Yao, Defu Lian, and Xing Xie. Recommender ai agent: Integrating large language models for interactive recommendations, 2024. URL https://arxiv.org/abs/2308.16505.
- [9] Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. Understanding the planning of llm agents: A survey, 2024. URL https://arxiv.org/abs/2402.02716.
- [10] Yanhong Li, Chenghao Yang, and Allyson Ettinger. When hindsight is not 20/20: Testing limits on reflective thinking in large language models, 2024. URL https://arxiv.org/abs/2404. 09129.
- [11] Bill Yuchen Lin, Yicheng Fu, Karina Yang, Prithviraj Ammanabrolu, Faeze Brahman, Shiyu Huang, Chandra Bhagavatula, Yejin Choi, and Xiang Ren. Swiftsage: A generative agent with fast and slow thinking for complex interactive tasks. *ArXiv preprint*, abs/2305.17390, 2023. URL https://arxiv.org/abs/2305.17390.
- [12] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=S37h0erQLB.
- [13] Keerthiram Murugesan, Mattia Atzeni, Pavan Kapanipathi, Pushkar Shukla, Sadhana Kumaravel, Gerald Tesauro, Kartik Talamadupula, Mrinmaya Sachan, and Murray Campbell. Text-based rl agents with commonsense knowledge: New challenges, environments and baselines. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(10):9018–9027, May 2021. doi: 10.1609/aaai.v35i10.17090. URL https://ojs.aaai.org/index.php/AAAI/ article/view/17090.
- [14] Karthik Narasimhan, Tejas Kulkarni, and Regina Barzilay. Language understanding for textbased games using deep reinforcement learning. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1–11, Lisbon, Portugal, September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1001. URL https: //aclanthology.org/D15-1001.
- [15] OpenAI. Gpt-4 technical report, 2023.
- [16] Matthew Renze and Erhan Guven. Self-reflection in llm agents: Effects on problem-solving performance, 2024. URL https://arxiv.org/abs/2405.06682.
- [17] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023.
- [18] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. ALFWorld: Aligning Text and Embodied Environments for Interactive Learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021. URL https://arxiv.org/abs/2010.03768.
- [19] Matthijs TJ Spaan. Partially observable markov decision processes. In *Reinforcement learning: State-of-the-art*, pages 387–414. Springer, 2012.

- [20] Ruo Yu Tao, Marc-Alexandre Côté, Xingdi Yuan, and Layla El Asri. Towards solving text-based games by producing adaptive action spaces, 2018.
- [21] Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. Science-World: Is your agent smarter than a 5th grader? In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11279–11298, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.775. URL https://aclanthology.org/2022.emnlp-main.775.
- [22] Zekun Wang, Ge Zhang, Kexin Yang, Ning Shi, Wangchunshu Zhou, Shaochun Hao, Guangzheng Xiong, Yizhi Li, Mong Yuan Sim, Xiuying Chen, Qingqing Zhu, Zhenzhu Yang, Adam Nik, Qi Liu, Chenghua Lin, Shi Wang, Ruibo Liu, Wenhu Chen, Ke Xu, Dayiheng Liu, Yike Guo, and Jie Fu. Interactive natural language processing, 2023. URL https://arxiv.org/abs/2305.13246.
- [23] Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents, 2024. URL https://arxiv.org/abs/2302.01560.
- [24] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. The rise and potential of large language model based agents: A survey, 2023. URL https://arxiv.org/abs/2309.07864.
- [25] Shunyu Yao, Rohan Rao, Matthew Hausknecht, and Karthik Narasimhan. Keep CALM and explore: Language models for action generation in text-based games. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8736–8754, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.704. URL https://aclanthology.org/2020.emnlp-main.704.
- [26] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models, 2023. URL https: //arxiv.org/abs/2210.03629.
- [27] Weiran Yao, Shelby Heinecke, Juan Carlos Niebles, Zhiwei Liu, Yihao Feng, Le Xue, Rithesh R N, Zeyuan Chen, Jianguo Zhang, Devansh Arpit, Ran Xu, Phil L Mui, Huan Wang, Caiming Xiong, and Silvio Savarese. Retroformer: Retrospective large language agents with policy gradient optimization. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=K0Zu91CzbK.
- [28] Ceyao Zhang, Kaijie Yang, Siyi Hu, Zihao Wang, Guanghe Li, Yihang Sun, Cheng Zhang, Zhaowei Zhang, Anji Liu, Song-Chun Zhu, Xiaojun Chang, Junge Zhang, Feng Yin, Yitao Liang, and Yaodong Yang. Proagent: Building proactive cooperative agents with large language models, 2024. URL https://arxiv.org/abs/2308.11339.
- [29] Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiuying Peng, Jun Wang, Yueting Zhuang, and Weiming Lu. Self-contrast: Better reflection through inconsistent solving perspectives, 2024. URL https://arxiv.org/abs/2401.02009.
- [30] Wenqi Zhang, Ke Tang, Hai Wu, Mengna Wang, Yongliang Shen, Guiyang Hou, Zeqi Tan, Peng Li, Yueting Zhuang, and Weiming Lu. Agent-pro: Learning to evolve via policy-level reflection and optimization, 2024. URL https://arxiv.org/abs/2402.17574.

## A Background

One may consider every TBG to be a partially observable Markov decision process (POMDP) [19] where the environment state is never observed directly. This POMDP may be formalized as  $\langle S, T, A, \Omega, R, \gamma \rangle$ , where  $\gamma \in [0, 1]$  denotes the reward discount factor. S denotes the set of states s that contain the internal information of the game – such as objects found throughout the game or the player's location – not all of which may be visible to the agent at any given time. A denotes the action space made up of individual text actions a issued by the player.  $\Omega$  denotes the observation function. Further,  $o \in O$  denotes the observations made by the player. The observation  $o_t$  of the agent at time t depends on the current state  $s_t$ , as well as the previous action  $a_{t-1}$ , which may be formalized as  $\Omega$  ( $o_t \mid s_t, a_{t-1}$ ). Seeing as the agent can only observe and interact with the environment of a TBG via natural language, each observation is composed of a sequence of tokens  $o_t = \{\hat{o}_t^1, \dots, \hat{o}_t^N\}$ , as are their actions  $a_t = \{\hat{a}_t^1, \dots, \hat{a}_t^M\}$ .

In the context of TBGs, an action  $a_t$  is considered admissible at a state  $s_t$  if it is capable of changing the game's state, *i.e.*, if it can lead to a transition to a new state  $s_{t+1}$  that is different from the current state  $s_t$ . The environment's state transition is modeled through a probabilistic function  $T(s_{t+1} | s_t, a_t)$ . Traditionally, admissible actions in state  $s_t$  could deterministically lead to a new state  $s_{t+1}$ . However, we use a more general approach where all actions, whether admissible or not, are included in the state transition function. Non-admissible actions, which do not lead to a change in the game's state, result in a transition back to the original state  $s_t$  with probability 1. In contrast, admissible actions lead to different states with their own probability. The admissible action set is bound to be significant for quests with a sufficiently large branching factor. While admissible, many action candidates are bound to be suboptimal.

The reward r received by the agent – the discounted sum of which,  $\mathbb{E}\left[\sum_{t} \gamma^{t} r_{t}\right]$ , it aims to maximize – are obtained by  $r_{t} = R(s_{t}, a_{t})$ . In practice, TBGs typically provide sparse positive rewards for completing subgoals and advancing towards completing the game.

A TBG's interpreter can accept any text sequence, but will only recognize text that follows a certain structure. Typically, commands take the form of either a single keyword (such as "look") or a combination involving verbs, objects, and occasionally prepositions. Previous works on TBGs made the assumption that we have access to a predefined set of all admissible actions at each game state and must select the correct action to progress [14, 6]. This is consistent with how some text adventure games are played in the real world [20]. In line with more recent works [11], we make the assumption that we have access to a number of action templates (*e.g.* connect A to B, pick up C) and subsequently ask the agent to generate the actions as tokens using these templates. This leads to a more challenging action generation process for the agent.

## **B** LLM Reflection

After successfully completing a sub task, the agent is prompted to reflect on its most recent observations and identify the key factors that contributed to its success. This reflection process encourages the agent to verbalize the steps and strategies that led to the desired outcome, helping it create a plan that can be stored in managed memory for future use. A similar procedure is applied in cases of failure. When a task is unsuccessful, the agent is instructed to reflect on alternative actions it could have taken and to devise a revised plan for the next attempt, ensuring continuous learning and improvement. This plan is also stored in memory.

An example trajectory of a ScienceWorld task where the agent reflects on positive and negative experiences is shown in Figure 2. At the end of this trajectory, we show for a single command what a reflection would look like for a successful or unsuccessful choice, in this case picking up an animal. For the unsuccessful case (red), where the agent does not pick up an animal, we reflect on whether another choice would have constituted an animal and thus resulted in a reward. For the successful case where a sub goal is reached (yellow), we reflect on what made the current actions successful and subsequently commit this to memory.

#### ScienceWorld Task 8-1

#### **Task Description:**

Your task is to find a(n) animal. First, focus on the thing. Then, move it to the red box in the kitchen.

> look around

This room is called the hallway. In it, you see:

- a substance called air
- a picture
- the agent

You also see:

- A door to the greenhouse (open)
- A door to the living room (open)
- A door to the art studio (open)
- A door to the kitchen (open)
- A door to the bedroom (open)
- A door to the workshop (open)

> open door to greenhouse

The door is already open.

> go to greenhouse

You move to the greenhouse.

> open door to outside

The door is already open.

> go to outside

You move to the outside.

> look around

This outside location is called the outside. Here, you see:

- a substance called air
- an axe
- a substance called wood
- the ground
- a fire pit (containing nothing)
- a fountain (containing a substance called water)
- a butterfly egg
- a blue jay egg
- a dove egg
- the agent

You also see:

- A door to the kitchen (open)
- A door to the foundry (open)
- A door to the greenhouse (open)

#### Case #1: Failure

> focus on axe

You focus on the axe.

> pick up axe

# You move the axe to the inventory.

Reflection(I failed because I assumed eggs are not animals. My plan was to go outside to find animals. I tried to pick up an axe. An alternative course of action would have been to go outside and pick up an egg, as it might be considered an animal.)

 $\rightarrow \text{Managed Memory}$ 

Case #2: Success > focus on dove egg You focus on the dove egg. > pick up dove egg You move the dove egg to the inventory. Reflection(Picking up an egg was the correct choice. Even though the animal was in an egg, it is still considered an animal. My plan was to go outside and find an animal in the wild. Items such as an axe, wood, or landscape elements are not related to animals.) → Managed Memory

Figure 2: Example of reflection in the 2-1 Task, where the agent needs to find something.

# C Example ScienceWorld Task

In this section, we provide a successfully completed task, a variation of Task 1-1, which concerns itself with boiling a substance, in this case water, to change its state. Once the agent has viewed its surroundings and moved to the kitchen, it collects the necessary items and begins its experiment by boiling the water in a pot on the stove. Finally, the agent examines steam and completes the task (highlighted in green). The example is truncated to improve readability. The trajectory is shown in Figure 3.

#### ScienceWorld Task 1-1

#### **Task Description:**

Your task is to boil water. For compounds without a boiling point, combusting the substance is also acceptable. First, focus on the substance. Then, take actions that will cause it to change its state of matter.

> look around

This room is called the hallway. In it, you see:

- a picture
- a substance called air
- the agent
- A door to the green house (open)
- A door to the living room (open)
- A door to the art studio (open)
- A door to the kitchen (open)
- A door to the bedroom (open)
- A door to the workshop (open)

> open door to kitchen

The door is already open.

> go to kitchen

You move to the kitchen. > look around

This room is called the kitchen. In it, you see:

- a substance called soap
- a painting
- a thermometer reading 10°C
- a counter with various items (banana, potato, red apple, orange)
- a sink (off, empty)

• a table with a glass cup (empty)							
• a chair (empty)							
• a freezer (closed)							
• a lighter							
- a righter							
• a stopwatch (deactivated)							
• a fridge (closed)							
• a cupboard (closed)							
• a stove (turned off)							
• a glass jar containing sodium chloride							
• various other items (air, the agent, etc.)							
<pre>&gt; pick up thermometer You move the thermometer to the inventory. &gt; open cupboard The cupboard is now open. &gt; pick up metal pot You move the metal pot to the inventory. &gt; move metal pot to sink You move the metal pot to the sink. &gt; activate sink The sink is now activated. &gt; deactivate sink The sink is now deactivated. &gt; pick up metal pot You move the metal pot to the inventory. &gt; focus on substance in metal pot You focus on the water. &gt; move metal pot to stove You move the metal pot to the stove. &gt; activate stove The stove is now activated.</pre>							
<ul> <li>&gt; use thermometer in inventory on substance in metal pot The thermometer measures a temperature of 13 degrees Celsius.</li> <li>&gt; examine substance in metal pot A substance called water.</li> </ul>							
<ul> <li>&gt; use thermometer in inventory on substance in metal The thermometer measures a temperature of 28 degrees Celsius.</li> <li>&gt; examine substance in metal pot A substance called water.</li> </ul>	pot						
<ul> <li>&gt; use thermometer in inventory on substance in metal The thermometer measures a temperature of 58 degrees Celsius.</li> <li>&gt; examine substance in metal pot A substance called water.</li> </ul>	pot						
<ul> <li>&gt; use thermometer in inventory on substance in metal The thermometer measures a temperature of 78 degrees Celsius.</li> <li>&gt; examine substance in metal pot A substance called water.</li> </ul>	pot						
<ul> <li>&gt; use thermometer in inventory on substance in metal pot The thermometer measures a temperature of 98 degrees Celsius.</li> <li>&gt; examine steam</li> </ul>							
A substance called steam. Task completed.							

Figure 3: Successful sequence of events in the 1-1 Task, where the agent needs to boil water.