# LOOKALIKE: HUMAN MIMICRY BASED COLLABORATIVE DECISION MAKING

### A PREPRINT

**◉ Rabimba Karanjai**\*
Department of Computer Science
University Of Houston
rkaranjai@uh.edu

**Weidong Shi**
Department of Computer Science
University Of Houston
wshi3@uh.edu

March 19, 2024

### ABSTRACT

Artificial General Intelligence falls short when communicating role specific nuances to other systems. This is more pronounced when building autonomous LLM agents capable and designed to communicate with each other for real world problem solving. Humans can communicate context and domain specific nuances along with knowledge, and that has led to refinement of skills. In this work we propose and evaluate a novel method that leads to knowledge distillation among LLM agents leading to realtime human role play preserving unique contexts without relying on any stored data or pretraining. We also evaluate how our system performs better in simulated real world tasks compared to state of the art.

*Keywords* LLM, AI, Machine Learning, Role Play, swarm, collaboration

## 1 Introduction

The rise of large, pre-trained models has changed the AI landscape. Transformer-based language models (LLMs) excel at natural language tasks. Researchers are now trying LLMs for sequential decision-making – a core skill for AI agents. While LLMs can suggest basic plans [Huang et al.(2022)], they often fail to consider real-world constraints or long-term planning. To improve this, approaches using environmental feedback are being explored. This feedback could be sensory data [Xiang et al.(2024)], human input [Yao et al.(2022)], or information about the plan's progress [Raman et al.(2022)], allowing the AI to adjust its plans accordingly.

LLMs show promise as planners, but they're far from perfect. A key hurdle is their reasoning ability; even given detailed instructions, they struggle to reliably produce functional plans [Silver et al.(2022)]. Current methods also force LLMs to learn through real-world (or simulated) trial and error, which is slow and costly. This differs from classical planning approaches [Fikes and Nilsson(1971)] where hypothetical plans can be analyzed for potential errors. Further, LLMs can be frustratingly stubborn, repeating the same mistakes despite feedback. Researchers are pushing the limits of what LLMs can learn and generalize, exploring new methods inspired by cognitive science to understand how information is shared within these models.

In this work, we explore different challenging reasoning scenarios involving both text based games and visual scenarios, more complex than previous studies. We explore a mimicking approach from the domain of cognitive science [Whiten(1998), Dawson and Foss(1965)] and see if it can be combined with more traditional imitation learning approaches [Torabi(2019)]. And through that we try to answer if we can introduce context aware problem solving skills to LLMs that are not capable of it without the framework.

---

\*www.rabimba.me

## 2  Problem Statement

The ability to gain knowledge, skills and solve relevant tasks and problems is generally defined as Intelligence[Chollet(2019)]. Such knowledge helps execute step by step tasks leading to a specific outcome. A lot of times these depend on a context of the task, to solve them effectively. Huamns gain these context and knowledge often from environment and other huamns, proving the context and knowledge transmissions as an effective tool in the context of problem solving. However in case of LLM's and LLM agents that is not the case always. The world knowledge, as available to us, is not available to a LLM making it much inefficient (and at times impossible) to solve certain problems. Knoweldge trasmission is a form of social learning, often assisted by other humans. This is common in everyday social interaction for humans. Corsss knoweldge transfers like copying a recipe by seeing someone cook, to learning how to play a game and learning rules of a game, are all cross domain knowledge transfer we do naturally. Which the LLMs fail to capture.

In our work, we seek to introduce our framework to capture the visual modal of these novel contexts to augment an LLMs understanding of a problem, and also cross transfer that to other agents preserving world knowledge for that task.

### 2.1  Motivation

The motivations for this work are the following:

- Increase domain specific reasoning capabilities of the LLM agents without data dependent training
- Rich visual scenes encode temporal and contextual knowledge that is beneficial and often crucial for solving a lot of real world tasks, which often is not encoded or captured by LLMs unelss specifically encoded.
- Contextual knowledge and knowledge transfer is a natural part of human intelligence. As such it is assumed they would also benefit AGI or to achieve AGI.

## 3  Reasoning Playground Design

### 3.1  Tasks

We specifically chose two tasks for our reasoning evaluation. One involving playing textual games based on world knowledge. Playing games has been one of the most evaluated tasks among the world planning scenarios [Tan et al.(2023)]. Artificial intelligence (AI) has made historic strides through games. A major milestone was IBM's Deep Blue defeating chess champion Garry Kasparov in 1997[Campbell et al.(2002)]. Google DeepMind's AlphaGo [Silver et al.(2016)] made history in 2016 by beating a professional Go player, a feat previously considered difficult for AI due to the game's complexity. AI continued to conquer poker with DeepStack[Moravčík et al.(2017)] and Libratus[Brown and Sandholm(2018)], programs that mastered heads-up no-limit Texas Hold'em in 2017. In 2019, OpenAI Five and DeepMind AlphaStar made waves by defeating world-class players in the complex strategy games Dota 2 [Berner et al.(2019)] and StarCraft II[Arulkumaran et al.(2019)], respectively. These victories underscore the increasing ability of AI to tackle increasingly challenging game environments.

To address the identified limitations [Sobieszek and Price(2022)] of the existing work . For our evaluation. we chose to use Zork [Hausknecht et al.(2020), Tsai et al.(2023)].

For a more comprehensive world planning we chose to utilize procgen [Cobbe et al.(2020)] to generate world scenarios for task and see if our farmework achieves goal in defined time.

## 4  Architecture

### 4.1  Learning

We believe the ability to faithfully learn from others is a powerful tool for LLMs – and we show that reinforcement learning can develop this skill solely from seeking rewards using a reward model actor. Our work requires just a few basic assumptions: that the environment includes a domain expert to mimic, and that the AI agent itself has some minimal ability to analyze and store information. The challenge is that the AI must figure out how the skilled individual acts without direct access to their decision-making process. We utilize a reward model[Silver et al.(2021)] mechanism to reward the AI agent to guide its actions towards favorable outcomes from the domain expert model. If our AI successfully learns to imitate, then the reinforcement learning process itself must have pushed for accurate, adaptable imitation that works in new situations. We demonstrate this using maximum a posteriori policy optimisation[Abdolmaleki et al.(2018)].
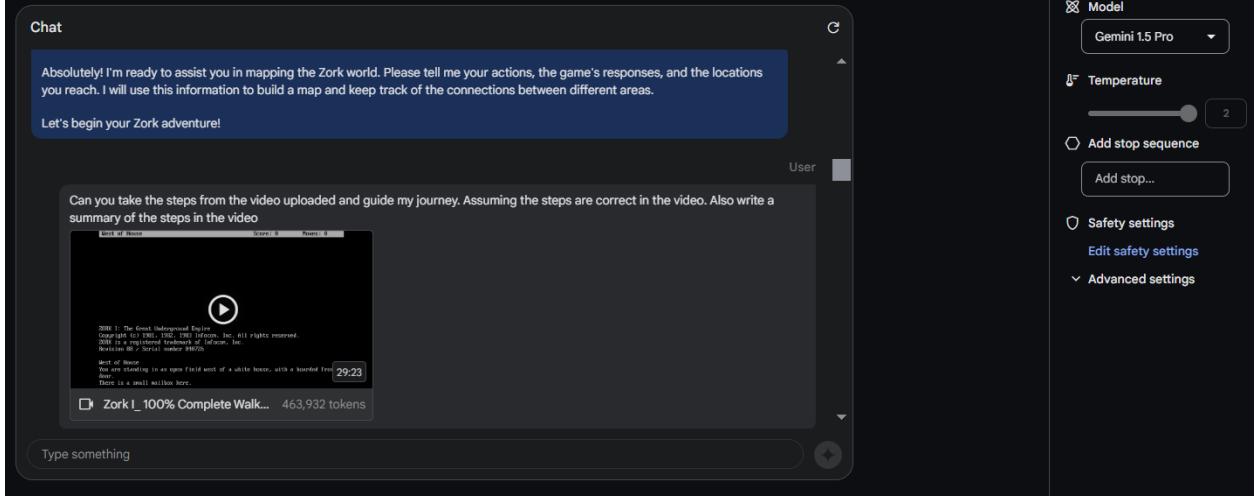
Figure 1: Zork Human Gameplay ingestion

In reinforcement learning (RL), there's a constant dillemma between trying new things (exploration) and using what you already know (exploitation). Exploration is vital for finding better strategies, but it can be risky, especially when rewards are rare or the problem is really complex[Kearns and Singh(2002)]. The risk is getting stuck doing something average instead of finding what works best. We demonstrate how an AI can learn to model its environment, including the actions of others, to tackle these difficult exploration challenges.

The trick is that standard approaches to RL don't work when you can't see everything that's going on or when the system is too complex to perfectly map out. In these cases, the AI needs to learn not only for itself but also by cleverly figuring out what the 'domain experts' around it are doing. We show that with the right setup, this kind of observational learning can happen through standard reinforcement learning, without any special instructions to imitate. This opens up new possibilities for solving truly challenging problems.

## 4.2 Playground

For Zork, we followed [Tsai et al.(2023)] closely to replicate the environment for playing the game. We took a similar approach to piping the question and answers from the game back to forth from the LLM. HOwever unlike their work, we have automated to process to keep make the experiment reproducible and run multiple occurrences.

For procgen we utilized their API to generate and solve the game models passing the task and task instructions through the LLM.

## 5    Evaluation Architecture

Our system comprises of tow sets of LLMs. One is the *player* LLM. (pLLM) which gets the environment information and comes up with the instructions to pass on to the game. This is true for both text based as well as vision based tasks. The player LLM (pLLM) in our case is gemma-7b[gem(2024)]. We intentionally chose a moderate sized LLM to see if our framework can actually utilize the encoded context information to get better apart from its own learned representations of knowledge.

The system also utilizes a separate Google Gemini 1.5 Pro model as a Domain Expert(vLLM). This model gets input from existing plays by human players to construct a reward mechanism for the pLLM. As an example, for zork, this gets input in few shot prompting of perfect Zork plays and extracts the gameplay information from the video as we see from Figure . Since Gemini 1.5 is capable of handling a very large context and also multimodal, this was perfectly suited for our automated domain expert critique.

Once it creates a outline of the steps, it can generate an optimal pathway for the game to play and can reward/punish the pLLM based on the existing reward. This vLLM (Vision LLM) is also used for vision related data ingestion for procgen tasks.

| GamePlay | # of actions | gemma-7b errors | gemini errors |
|----------|--------------|-----------------|---------------|
| Map Navigation | 27 | 248+ | 17 |
| Decisive Action | 15 | 88+ | 6 |

Table 1: Table 1: The number of errors in the gameplay steps produced by the LLMs. A "+" mark means the generated model contains too much irrelevant information, which obscures the true number of errors.

### 5.1 Constructing Game Play Solution

If we generalize our task, our solution starts by feeding pLLM a variety of information. We start with detailed instructions on how to play the game. Example videos from similar tasks are fed to vLLM for building a reward model. Next, we provide context about the current situation, including the agent's goals and any game limitations. This is followed by a description of the specific action the agent is taking. Finally, we include a list of reusable terms so the LLM maintains a consistent language within the plan. Each of the steps is rewarded/punished based on the outcome as determined by the vLLM model.

An AI agent must first acknowledge the existence of other players [Ndousse et al.(2020)] and the likelihood that they share its objectives and capabilities in order to use information from them successfully. The agent won't, however, have direct access to the thoughts or acts of the other players. In order to solve this, we give the AI a unique "attention" mechanism that aids in keeping it focused on the locations of other players in its surroundings. Similar to a spotlight, this attention method aids the AI in recognizing what's crucial, in this case, maintaining track of the other participants. When the AI has mastered paying attention to the appropriate things, we no longer need the unique sensor that we used during training to assist it learn this skill. This is only done for pLLM, with different temperature for it to assume different role. Once pLLMs with different temperature goes through the rewards of vLLM, we take into account the most successful one.

In case of pLLM swarms, for each step, after each reward/punishment for each of the pLLMs for each step, the results are stored and passed back as "memory" to the pLLMs along with the rewards/ punishment information for the next step, before being rewarded.

## 6 Results

We tested our approach in three settings: just asking a pLLM on the game play, and then on two settings, with vLLM reward and one with both vLLM reward as well as swarm input. POur swarm evaluation setup is more complex than other benchmarks, with the pLLM having more inputs to consider as piori than just a reward being given *after* the decision.

Our experiments focused on a few key areas. First, we looked at how well the LLMs generated gameplay is. Next, we tested if the LLMs could use collaborative feedback from piror steps to fix errors in their play. Since the pLLMs were run in different temperature (and intentionally not trying to be reproducible in their respective plays) the feedback from the reward models varied. We looked at if these were beneficial for another pLLM to imporve their game play assuming a different role. Finally, we demonstrated how the just giving vLLM produced rewards along with pLLM swarm could be used by the AI for actual planning and problem-solving. We used two language models for this: gemma-7b and Gemini 1.5 Pro.

### 6.1 Plan Generation

We started by testing how well the language models could come up with game plan when only some of the rules were provided. This is important because in real-world situations, you often know some constraints about the possible actions beforehand. We focused on whether the models could accurately recreate a "correct" game plan that took those constraints and the relationships between actions into account. Our results show that without reward Gemini significantly outperforms gemma-7b in generating accurate plans. Table 1 details the number of errors for each domain, and to give you an idea of the complexity involved, we've also included the number of parameters and literals in the final, corrected Gemini plans. We can see the errors in Table 1.

| Player Type | Zork | procgen |
|---|---|---|
| pLLM | 35 | 0 |
| pLLM+vLLM | 37 | 15 |
| pLLM Swarm + vLLM | 100 | 56 |

Table 2: Table : Success rates of different approaches in solving a gameplay

### 6.2 Rewarding using vLLM

We decided to focus on the reward models Gemini produced with only videos as input. Our goal was to show that Gemini could be used as a tool to take video and fix errors in its plans mimicking the human gameplay. These mimicking behaviours would act as reward mechanism for pLLMs swarm decision making collaboration.

By providing feedback like "Possible actions: Open mailbox Go north Go south Go west Recommended action: Open the mailbox." Gemini could pinpoint and define decisive actions. These action rewards successfully corrected a lot of incorrect outputs from the pLLMs and that in turn helped others in the swarm to reach the solution faster.

For different tasks we can see in Table 2 our method remarkably improves plan success rate, and solves Zork comfortably. Even for procgen tasks, which the standalone pLLM couldn't solve at all, show an increased success.

## 7 Discussion

There are two perspectives on open-endedness in LLM systems for problem solving. We can train an LLM to achieve a specific degree of success over a wide range of tasks that get harder and harder, or we can concentrate on one complex activity and aim to create an LLM that keeps getting better at it. For our use case—coming up with solutions for game-driven tasks—we found the second technique to be more useful.

Our framework itself is surprisingly basic. This shows that the way to create flexible, "collaborative" LLMs that can correct themselves from cooperative interactions may be to use fewer, more potent components. In the future, it will be essential to develop scalable training techniques, enhance how AIs learn about the world, and build systems in which the training process itself is dynamic.

Our method of using social learning to teach the LLM is similar to "memory-based" super-learning. As an example in the event that a new book in a library or if the catalogue shifts, you would want a librarian robot to be able to rapidly adapt. Sending all of its data to a massive server may also raise privacy issues. Our LLM absorbs human knowledge by emulating human behavior and storing all of it in its own memory. This maintains privacy while making it flexible.

We have demonstrated that social learning research can be conducted in complex visual contexts as well. But scale alone isn't the key. A scaled-down version of our setup could also provide exciting findings. The scalability of the approach is a powerful incentive for further research in this field, regardless of one's access to massive computing resources.

## 8 Conclusion

In this work we have introduced a new way of using language models (LLMs) for planning and solving tasks. Instead of having them generate plans directly, we use LLMs to build a model of how the world works by participating in a collaborative swarm and rewarding correct actions based on na domain expert(vLLM). This approach leverages the strength of LLMs – understanding the world – while avoiding their weakness in complex planning. Our process starts with Gemini generating a detailed vision guided reward plan. We then use this to refine the pLLMs. Finally, use this collaborative decision to play and evaluate the games. Each of these pLLMs assume a different role for a gamer, and achieves efficiency by collaborating and passing contextual information.

There are rooms for improvement in our system. First, our test environments were still relatively simple compared to some planning challenges. Can LLMs scale to more complex logic? Second, we currently assume the the videos are

enough for domain expert and procgen correctly generates enough variation for us to ignore requirement for Automatic Domain Randomisation in our work. However a more detailed work with ADR is required to prove generalisation of knowledge transfer between the pLLM agents. Finally, we assume perfect knowledge of object states, but the real world has messy perception, and our system needs to account for that.

# References

[gem(2024)] 2024. *Gemma: Google introduces new state-of-the-art open models.* `https://blog.google/technology/developers/gemma-open-models/`

[Abdolmaleki et al.(2018)] Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Remi Munos, Nicolas Heess, and Martin Riedmiller. 2018. Maximum a posteriori policy optimisation. *arXiv preprint arXiv:1806.06920* (2018).

[Arulkumaran et al.(2019)] Kai Arulkumaran, Antoine Cully, and Julian Togelius. 2019. Alphastar: An evolutionary computation perspective. In *Proceedings of the genetic and evolutionary computation conference companion.* 314–315.

[Berner et al.(2019)] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. 2019. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680* (2019).

[Brown and Sandholm(2018)] Noam Brown and Tuomas Sandholm. 2018. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science* 359, 6374 (2018), 418–424.

[Campbell et al.(2002)] Murray Campbell, A Joseph Hoane Jr, and Feng-hsiung Hsu. 2002. Deep blue. *Artificial intelligence* 134, 1-2 (2002), 57–83.

[Chollet(2019)] François Chollet. 2019. On the measure of intelligence. *arXiv preprint arXiv:1911.01547* (2019).

[Cobbe et al.(2020)] Karl Cobbe, Chris Hesse, Jacob Hilton, and John Schulman. 2020. Leveraging procedural generation to benchmark reinforcement learning. In *International conference on machine learning.* PMLR, 2048–2056.

[Dawson and Foss(1965)] Betty V Dawson and BM Foss. 1965. Observational learning in budgerigars. *Animal behaviour* (1965).

[Fikes and Nilsson(1971)] Richard E Fikes and Nils J Nilsson. 1971. STRIPS: A new approach to the application of theorem proving to problem solving. *Artificial intelligence* 2, 3-4 (1971), 189–208.

[Hausknecht et al.(2020)] Matthew Hausknecht, Prithviraj Ammanabrolu, Marc-Alexandre Côté, and Xingdi Yuan. 2020. Interactive fiction games: A colossal adventure. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 7903–7910.

[Huang et al.(2022)] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. 2022. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning.* PMLR, 9118–9147.

[Kearns and Singh(2002)] Michael Kearns and Satinder Singh. 2002. Near-optimal reinforcement learning in polynomial time. *Machine learning* 49 (2002), 209–232.

[Moravčík et al.(2017)] Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisỳ, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. 2017. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science* 356, 6337 (2017), 508–513.

[Ndousse et al.(2020)] Kamal Ndousse, Douglas Eck, Sergey Levine, and Natasha Jaques. 2020. Multi-agent social reinforcement learning improves generalization. *arXiv preprint arXiv:2010.00581* (2020).

[Raman et al.(2022)] Shreyas Sundara Raman, Vanya Cohen, Eric Rosen, Ifrah Idrees, David Paulius, and Stefanie Tellex. 2022. Planning with large language models via corrective re-prompting. In *NeurIPS 2022 Foundation Models for Decision Making Workshop.*

[Silver et al.(2016)] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. 2016. Mastering the game of Go with deep neural networks and tree search. *nature* 529, 7587 (2016), 484–489.

[Silver et al.(2021)] David Silver, Satinder Singh, Doina Precup, and Richard S Sutton. 2021. Reward is enough. *Artificial Intelligence* 299 (2021), 103535.

[Silver et al.(2022)] Tom Silver, Varun Hariprasad, Reece S Shuttleworth, Nishanth Kumar, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. 2022. PDDL planning with pretrained large language models. In *NeurIPS 2022 Foundation Models for Decision Making Workshop*.

[Sobieszek and Price(2022)] Adam Sobieszek and Tadeusz Price. 2022. Playing games with AIs: the limits of GPT-3 and similar large language models. *Minds and Machines* 32, 2 (2022), 341–364.

[Tan et al.(2023)] Qinyue Tan, Ashkan Kazemi, and Rada Mihalcea. 2023. Text-Based Games as a Challenging Benchmark for Large Language Models. (2023).

[Torabi(2019)] Faraz Torabi. 2019. Imitation learning from observation. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence* (Honolulu, Hawaii, USA) *(AAAI'19/IAAI'19/EAAI'19)*. AAAI Press, Article 1249, 2 pages. `https://doi.org/10.1609/aaai.v33i01.33019900`

[Tsai et al.(2023)] Chen Feng Tsai, Xiaochen Zhou, Sierra S Liu, Jing Li, Mo Yu, and Hongyuan Mei. 2023. Can Large Language Models Play Text Games Well? Current State-of-the-Art and Open Questions. *arXiv preprint arXiv:2304.02868* (2023).

[Whiten(1998)] Andrew Whiten. 1998. Imitation of the sequential structure of actions by chimpanzees (Pan troglodytes). *Journal of comparative psychology* 112, 3 (1998), 270.

[Xiang et al.(2024)] Jiannan Xiang, Tianhua Tao, Yi Gu, Tianmin Shu, Zirui Wang, Zichao Yang, and Zhiting Hu. 2024. Language models meet world models: Embodied experiences enhance language models. *Advances in neural information processing systems* 36 (2024).

[Yao et al.(2022)] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629* (2022).