# POKÉLLMON: A Human-Parity Agent for Pokémon Battles with Large Language Models

**Sihao Hu, Tiansheng Huang, Ling Liu**

Georgia Institute of Technology

Atlanta, GA 30332, United States

{sihaohu, thuang, ling.liu}@gatech.edu

https://poke-llm-on.github.io/

## Abstract

We introduce POKÉLLMON, the first LLM-based agent that achieves human-parity performance in tactical battle games, as demonstrated in Pokémon battles. The design of POKÉLLMON incorporates three key strategies: (i) In-context reinforcement learning that instantly consumes text-based feedback derived from battles to iteratively refine the policy; (ii) Knowledge-augmented generation that retrieves external knowledge to counteract hallucination and enables the agent to act timely and properly; (iii) Consistent action generation to mitigate the *panic switching* phenomenon when the agent faces a powerful opponent and wants to elude the battle. We show that online battles against human demonstrates POKÉLLMON's human-like battle strategies and just-in-time decision making, achieving 49% of win rate in the Ladder competitions and 56% of win rate in the invited battles. Our implementation and playable battle logs are available at: https://github.com/git-disl/PokeLLMon.

## 1. Introduction

Generative AI and Large Language Models (LLMs) have shown unprecedented success on NLP tasks (Ouyang et al., 2022; Brown et al., 2020; Xi et al., 2023; Wang et al., 2023b). One of the forthcoming advancements will be to explore how LLMs can autonomously act in the physical world with extended generation space from text to action, representing a pivotal paradigm in the pursuit of Artificial General Intelligence (Goertzel & Pennachin, 2007; Goertzel, 2014).

Games are suitable test-beds to develop LLM-based agents (Duan et al., 2022; Batra et al., 2020) to interact with the virtual environment in a way resembling human behavior. For example, Generative Agents (Park et al., 2023) conducts a social experiments with LLMs assuming various
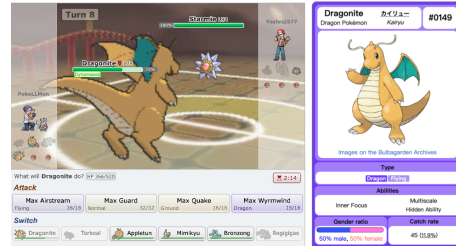


Figure 1. At each turn, the player is requested to decide which action to perform, *i.e.*, whether to let *Dragonite* to take a move or switch to another Pokémon off the field.

roles in a "The Sims"-like sandbox, where agents exhibit behavior and social interactions mirroring humans. In Minecraft, decision-making agents (Wang et al., 2023a;c; Singh et al., 2023) are designed to explore the world and develop new skills for solving tasks and making tools.

Compared to existing games, tactical battle games (Ma et al., 2023) are better suited for benchmarking the game-playing ability of LLMs as the win rate can be directly measured and consistent opponents like AI or human players are always available. Pokémon battles, serving as a mechanism that evaluates the battle abilities of trainers in the well-known Pokémon games, offer several unique advantages as the first attempt for LLMs to play tactical battle games:

(1) The state and action spaces are discrete and can be translated into text losslessly. **Figure** 1 is an illustrative example for a Pokémon battle: At each turn, the player is requested to generate an action to perform given the current state of Pokémon from each side. The action space consists of four moves and five possible Pokémon to switch; (2) The turn-based format eliminates the demands of intensive gameplay, alleviating the stress on the inference time cost for LLMs, making performance hinges solely on the reasoning abilities of LLMs; (3) Despite its seemingly simple mechanism, Pokémon battle is strategic and complex: an experienced player takes various factors into consideration, including species/type/ability/stats/item/moves of all the Pokémon on and off the field. In a random battle, each Pokémon is

randomly selected from a large candidate pool (more than 1,000) with distinct characteristics, demanding the players both the Pokémon knowledge and reasoning ability.

**Scope and Contributions:** The scope of this paper is to develop an LLM-based agent that mimics the way a human player engages in Pokémon battles. The objective is to explore the key factors that make the LLM-based agent a good player and to examine its strengths and weaknesses in battles against human players. To enable LLMs play game autonomously, we implement an environment that can parse and translate battle state into text description, and deliver generated action back to the server. By evaluating existing LLMs, we identify the presence of *hallucination*, and the *panic switching* phenomenon.

*Hallucination*: The agent can mistakenly send out Pokémon at a type disadvantage or persist in using ineffective moves against the opponent. As a result, the most advanced LLM, GPT-4, achieves a win rate of 26% when playing against a heuristic bot, compared to 60% win rate of human players. To combat hallucination, we introduce two strategies: (1) In-context reinforcement learning: We provide the agent with text-based feedback *instantly* derived from the battle, serving as a new form of "reward" to *iteratively* refine the action generation policy without training; (2) Knowledge-augmented generation: We equip the agent with Pokédex, an encyclopaedia in Pokémon games that provides external knowledge like type advantage relationship or move/ability descriptions, simulating a human player searching for the information of unfamiliar Pokémon.

*Panic switching*: We discover that when the agent encounters a powerful Pokémon, it tends to panic and generates inconsistent actions like switching different Pokémon in consecutive turns to elude the battle, a phenomenon that is especially pronounced with Chain-of-Thought ([Wei et al., 2022](#)) reasoning. Consistent action generation alleviates the issue by voting out the most consistent action without overthinking. This observation mirrors human behavior, where in stressful situations, overthinking and exaggerating difficulties can lead to panic and impede acting.

Online battles demonstrate POKÉLLMON's human-competitive battle abilities: it achieves a 49% win rate in the Ladder competitions and a 56% win rate in the invited battles. Furthermore, we reveal its vulnerabilities to human players' attrition strategies and deceptive tricks.

In summary, this paper makes four original contributions:

- We implement and release an environment that enables LLMs to autonomously play Pokémon battles.

- We propose in-context reinforcement learning to instantly and iteratively refine the policy, and knowledge-augmented generation to combat hallucination.

- We discover that the agent with chain-of-thought experiences panic when facing powerful opponents, and consistent action generation can mitigate this issue.

- POKÉLLMON, to the best of our knowledge, is the first LLM-based agent with human-parity performance in tactical battle games.

## 2. LLMs as Game Players

**Communicative games:** Communicative games revolve around communication, deduction and sometimes deception between players. Existing studies show that LLMs demonstrate strategic behaviors in board games like Werewolf ([Xu et al., 2023](#)), Avalane ([Light et al., 2023](#)), WorldWar II ([Hua et al., 2023](#)) and Diplomacy ([Bakhtin et al., 2022](#)).

**Open-ended games:** Open-ended games allow players to freely explore the game world and interact with others. Generative Agent ([Park et al., 2023](#)) showcases that LLM-based agents exhibit behavior and social interactions mirroring human-like patterns. In MineCraft, Voyager ([Wang et al., 2023a](#)) employs curriculum mechanism to explore the world and generates and executes code for solving tasks. DEPS ([Wang et al., 2023c](#)) proposes an approach of "Describe, Explain, Plan and Select" to accomplish 70+ tasks. Planing-based frameworks like AutoGPT ([Significant Gravitas](#)) and MetaGPT ([Hong et al., 2023](#)) can be adopted for the exploration task as well.

**Tactic battle games:** Among various game types, tactical battle games ([Akata et al., 2023](#); [Ma et al., 2023](#)) are particularly suitable for benchmarking LLMs' game-playing ability, as the win rate can be directly measured, and consistent opponents are always available. Recently, LLMs are employed to play StarCraft II ([Ma et al., 2023](#)) against the built-in AI with a text-based interface and a chain-of-summarization approach. In comparison, POKÉLLMON has several advantages: (1) Translating Pokémon battle state into text is lossless; (2) Turn-based format eliminates real-time stress given the inference time cost of LLMs; (3) Battling against disciplined human players elevates the difficulty to a new height.

## 3. Background

### 3.1. Pokémon

**Species:** There are more than 1,000 Pokémon species ([bul, 2024c](#)), each with its unique ability, type(s), statistics (stats) and battle moves. **Figure** [2](#) shows two representative Pokémon: *Charizard* and *Venusaur*.

**Type:** Each Pokémon species has up to two elemental types, which determine its advantages and weaknesses. **Figure** [3](#) shows the advantage/weakness relationship between 18 types of attack moves and attacked Pokémon. For ex-

*Figure 2.* Two representative Pokémon: *Charizard* and *Venusaur*. Each Pokémon has type(s), ability, stats and four battle moves.

ample, fire-type moves like "Fire Blast" of *Charizard* can cause double damage to grass-type Pokémon like *Venusaur*, while *Charizard* is vulnerable to water-type moves.

**Stats:** Stats determine how well a Pokémon performs in battles. There are four stats: (1) Hit Points (HP): determines the damage a Pokémon can take before fainting; (2) Attack (Atk): affects the strength of attack moves; (3) Defense (Def): dictates resistance against attacks; (4) Speed (Spe): determines the order of moves in battle.

**Ability:** Abilities are passive effects that can affect battles. For example, *Charizard*'s ability is "Blaze", which enhances the power of its fire-type moves when its HP is low.

**Move:** A Pokémon can learn four battle moves, categorized as attack moves or status moves. An attack move deals instant damage with a power value and accuracy, and associated with a specific type, which often correlates with the Pokémon's type but does not necessarily align; A status move does not cause instant damage but affects the battle in various ways, such as altering stats, healing or protect Pokémon, or battle conditions, *etc*. There are 919 moves in total with distinctive effect (bul, 2024b).

### 3.2. Battle Rule

In one-to-one random battles (Wikipedia, 2023), two battlers face off, each with six randomly selected Pokémon. Initially, each battler sends out one Pokémon onto the field, keeping the others in reserve for future switches. The objective is to make all the opponent's Pokémon faint (by reducing their HP to zero) while ensuring that at least one of own Pokémon remains unfainted. The battle is turn-based: at the start of each turn, both players choose an action to perform. Actions fall into two categories: (1) taking a move, or (2) switching to another Pokémon. The battle engine executes actions and updates the battle state for the next step. If a Pokémon faints after a turn and the battler has other Pokémon unfainted, the battle engine forces a switch, which does not count as the player's action for the next step. After a forced switch, the player can still choose a move or make another switch.
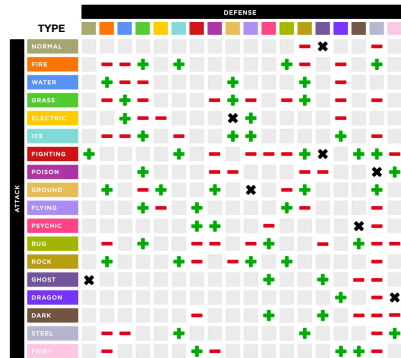


*Figure 3.* Type advantage/weakness relationship. "+" denotes super-effective (2x damage); "−" denotes ineffective (0.5x damage); "×" denotes no effect (0x damage). Unmarked is standard (1x) damage.

## 4. Battle Environment

**Battle Engine:** The environment interacts with a battle engine server called Pokémon showdown (pok, 2024), which provides a web-based GUI for human players, as well as web APIs for interacting with message in defined formats.

**Battle Environment:** We implement a battle environment based on (Sahovic, 2023a) to support LLMs autonomously play Pokémon battles. **Figure** 4 illustrates how the entire framework works. At the beginning of a turn, the environment get an action-request message from the server, including the execution result from the last turn. The environment first parses the message and update local state variables, and then translates the state variables into text. The text description primarily consists of four parts: (1) Own team information, including the attributes of Pokémon both on-the-field and off-the-field; (2) Opponent team information including the attributes of opposing Pokémon on-the-field and off-the-field (some are unknown); (3) Battle field information like the weather, entry hazard and terrain; (4) Historical turn log information, including previous actions of both side Pokémon, which is stored in a log queue. LLMs take the translated state as input and output an action for the next step. The action is sent to the server and executed alongside the action chosen by the human player.

## 5. Preliminary Evaluation

To gain insights into the challenges associated with Pokémon battles, we evaluate the abilities of existing LLMs, including GPT-3.5 (Ouyang et al., 2022), GPT-4 (Achiam et al., 2023), and LLaMA-2 (Touvron et al., 2023),

### 5.1. Pokémon Battles

Placing LLMs in direct competitions against human players is time-consuming as human needs time to think (4 minutes for 1 battle in average). To save time, we adopt a heuris-
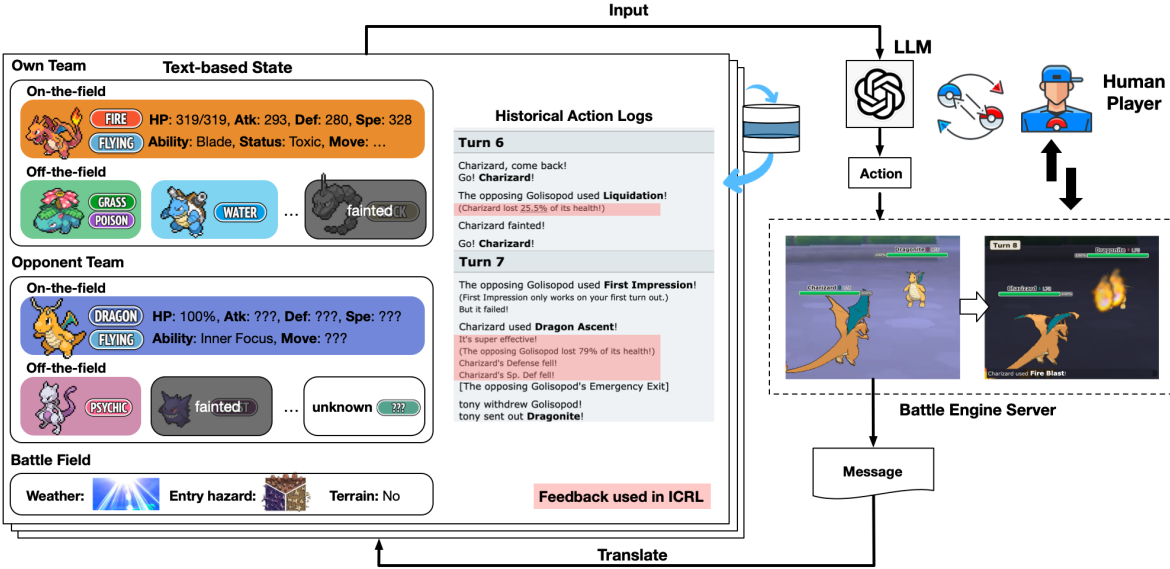
*Figure 4.* The framework that enables LLMs to battle with human players: It parses the messages received from the battle server and translates state logs into text. LLMs take these state descriptions and historical turn logs as input and generates an action for the next step. The action is then sent to the battle server and executed alongside the action chosen by the opponent player.

tic bot (Sahovic, 2023b) to initially battle against human players in the Ladder competitions, and then use the bot to benchmark existing LLMs. The bot is programmed to use status boosting moves, set entry hazards, selecting the most effective actions by considering the stats of Pokémon, the power of moves, and type advantages/weaknesses.

*Table 1.* Performance of LLMs in battles against the bot.

| Player | Win rate ↑ | Score ↑ | Turn # | Battle # |
|---|---|---|---|---|
| Human | 59.84% | 6.75 | 18.74 | 254 |
| Random | 1.2% | 2.34 | 22.37 | 200 |
| MaxPower | 10.40% | 3.79 | 18.11 | 200 |
| LLaMA-2 | 8.00% | 3.47 | 20.98 | 200 |
| GPT-3.5 | 4.00% | 2.61 | 20.09 | 100 |
| GPT-4 | 26.00% | 4.65 | 19.46 | 100 |

The statistic results are presented in **Table 1**, where the battle score is defined as the sum of the numbers of the opponent's fainted Pokémon and the player's unfainted Pokémon at the end of a battle. Consequently, the opponent player's battle score is equal to 12 minus the player's battle score. Random is a simple strategy that randomly generates an action every time, and MaxPower chooses the move with the highest power value. Obviously, GPT-3.5 and LLaMA-2 are just slightly better than Random and even GPT-4 cannot beat the bot, let along well-disciplined human players from the Ladder competitions.

By observing LLMs play battles and analyzing the explanations generated with their actions, we identify the occurrence of hallucination (Rawte et al., 2023; Cabello et al., 2023): LLMs can mistakenly claim non-existent type-advantage relationships or, even worse, reverse the advan-

tage relationships between types like sending a grass-type Pokémon to face with a fire-type Pokémon. A clear understanding of type advantage/weakness is crucial in Pokémon battles, as choosing a Pokémon with a type advantage can result in dealing more damage and sustaining less.

## 5.2. Test of Hallucination

To assess hallucination in the outputs of LLMs, we construct the task of type advantage/weakness prediction. The task involves asking LLMs to determine if an attack of a certain type is A. super-effective (2x damage), B. standard (1x damage), C. ineffective (0.5x damage) or D. no effect (0x damage) against a certain type of Pokémon. The 324 (18x18) testing pairs are constructed based on **Figure 3**.

*Table 2.* Confusion matrices for type advantage prediction.

| Model | LLaMA-2 | | | | GPT-3.5 | | | | GPT-4 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Class | A | B | C | D | A | B | C | D | A | B | C | D |
| A | 5 | 46 | 0 | 0 | 0 | 0 | 49 | 2 | 37 | 8 | 5 | 1 |
| B | 25 | 179 | 0 | 0 | 2 | 6 | 185 | 11 | 0 | 185 | 17 | 2 |
| C | 15 | 46 | 0 | 0 | 0 | 2 | 57 | 2 | 3 | 24 | 32 | 2 |
| D | 1 | 7 | 0 | 0 | 0 | 0 | 7 | 1 | 0 | 0 | 0 | 8 |

**Table 2** shows the three confusion matrices of LLMs, where their performance is highly related to their win rates in **Table 1**. LLaMA-2 and GPT-3.5 suffer from severe hallucination problems, while GPT-4 achieves the best performance with an accuracy of 84.0%, we still observe it frequently making ineffective actions, which is because in a single battle, LLMs need to compare the types of all the opponent's Pokémon with types of all their Pokémon, as well as types of moves.
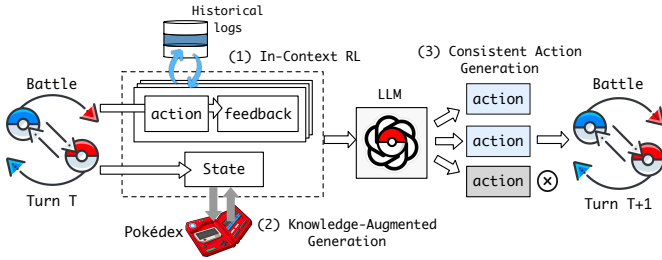
Figure 5. POKÉLLMON is equipped with three strategies: (1) ICRL that leverages instant feedbacks from the battle to iteratively refine generation; (2) KAG that retrieves external knowledge to combat hallucination and to act timely and properly; (3) Consistent Action Generation to prevent the panic switching problem.

# 6. POKÉLLMON

**Overview:** The overall framework of POKÉLLMON is illustrated in **Figure** 5. In each turn, POKÉLLMON uses previous actions and corresponding text-based feedback to *iteratively* refine the policy, and also augments the current state information with external knowledge, such as type advantage/weakness relationships and move/ability effects. Given above information as input, it independently generates multiple actions and selects the most consistent ones as the final output for execution.

## 6.1. In-Context Reinforcement Learning (ICRL)

Human players make decisions based not only on the current state but also on the (implicit) feedback from previous actions, such as the change in a Pokémon's HP over two consecutive turns following an attack by a move. Without feedback provided, the agent can continuously stick to the same erroneous action. As illustrated in **Figure** 6, the agent uses "Crabhammer", a water-type attack move against the opposing *Toxicroak*, a Pokémon with the ability "Dry Skin", which can nullify damage from water-type moves. The "Immune" message displayed in the battle animation can prompt a human player to change actions even without knowledge of "Dry Skin", however, is not included in the state description. As a result, the agent repeats the same action, inadvertently giving the opponent two free turns to triple *Toxicroak*'s attack stats, leading to defeat.

Reinforcement Learning (Schulman et al., 2017; Mnih et al., 2016; Hafner et al., 2023) requires numeric rewards to evaluate actions for refining policy. As LLMs can understand languages and distinguish what is good and bad, text-based feedback description provides a new form of "reward". By incorporating text-based feedback from the previous turns into the context, the agent is able to refine its "policy" *iteratively* and *instantly* during serving, namely In-Context Reinforcement Learning (ICRL).

In practice, we generate four types of feedback: (1) The



Figure 6. The agent repeatedly uses the same attack move but has zero effect to the opposing Pokémon due to its ability "Dry Skin."



Figure 7. In turn 3, the agent uses "Psyshock", which cause zero damage to the opposing Pokémon. With ICRL, the agent switch to another Pokémon.

change in HP over two consecutive turns, which reflects the actual damage caused by an attack move; (2) The effectiveness of attack moves, indicating whether they are super-effective, ineffective, or have no effect (immunity) due to type advantages or ability/move effects; (3) The priority of move execution, providing a rough estimate of speed, as precise stats for the opposing Pokémon are unavailable; (4) The actual effects of executed moves: both status and certain attack moves can cause outcomes like stat boosts or debuffs, recover HP, inflict conditions such as poison, burns, or freezing, *etc*. **Figure** 4 presents several instances of generated text-based feedback for ICLR.

Table 3. Performance of ICRL in battles against the bot.

| Player | Win rate ↑ | Score ↑ | Turn # | Battle # |
|--------|-----------|---------|--------|----------|
| Human  | 59.84%    | 6.75    | 18.74  | 254      |
| Origin | 26.00%    | 4.65    | 19.46  | 100      |
| ICRL   | 36.00%    | 5.25    | 20.64  | 100      |

**Table** 3 shows the improvement brought by ICRL. Compared to the original performance of GPT-4, the win rate is boosted by 10%, and the battle score increases by 12.9%. During the battles, we observe that the agent begins to change its action if the moves in previous turns do not meet the expectation, as shown in **Figure** 7: After observing that the opposing Pokémon is immune to the attack, it switches to another Pokémon.

## 6.2. Knowledge-Augmented Generation (KAG)

Although ICRL can mitigate the impact of hallucination, it can still cause fatal consequences before the feedback is received. For example, if the agent sends out a grass-type Pokémon against a fire-type Pokémon, the former is likely be defeated in a single turn before the agent realize it is a bad decision. To further reduce hallucination, Retrieval-Augmented Generation (Lewis et al., 2020; Guu et al., 2020; Patil et al., 2023) employ external knowledge to augment

5

_Table 4._ Performance of KAG in battles against the bot.

| Player | Win rate ↑ | Score ↑ | Turn # | Battle # |
|---|---|---|---|---|
| Human | 59.84% | 6.75 | 18.74 | 254 |
| Origin | 36.00% | 5.25 | 20.64 | 100 |
| KAG[Type] | 55.00% | 6.09 | 19.28 | 100 |
| KAG[Effect] | 40.00% | 5.64 | 20.73 | 100 |
| KAG | 58.00% | 6.53 | 18.84 | 100 |



_Figure 8._ The agent understands the move effect and uses it properly: _Klefki_ is vulnerable to the ground-type attack of _Rhydon_. Instead of switching, the agent uses "Magnet Rise", a move that protects itself from the ground-type attack for five turns, invalidating the ground-type attack "Earthquake" of the opposing _Rhydon_.

generation. In this section, we introduce two types of external knowledge to fundamentally mitigate hallucination.

**Type advantage/weakness relationship:** In the original state description in **Figure** 4, we annotate all the type information of Pokémon and moves to let the agent infer the type advantage relationship by itself. To reduce the hallucination contained in the reasoning, we explicitly annotate the type advantage and weakness of the opposing Pokémon and our Pokémon with descriptions like "_Charizard_ is strong against grass-type Pokémon yet weak to the fire-type moves".

**Move/ability effect:** Given the numerous moves and abilities with distinct effects, it is challenging to memorize all of them even for experienced human players. For instance, it's difficult to infer the effect of a status move based solely on its name: "Dragon Dance" can boost the user's attack and speed by one stage, whereas "Haze" can reset the boosted stats of both Pokémon and remove abnormal statuses like being burnt. Even attack moves can have additional effects besides dealing damage.

We collect all the effect descriptions of moves, abilities from Bulbapedia (bul, 2024b;a) and store them into a Pokédex, an encyclopaedia in Pokémon games. For each Pokémon on the battlefield, its ability effect and move effects are retrieved from the Pokédex and added to the state description.

**Table** 4 shows the results of generations augmented with two types of knowledge, where type advantage relationship (KAG[Type]) significantly boosts the win rate from 36% to 55%, whereas, Move/ability effect descriptions also enhance the win rate by 4 AP. By combining two of them, KAG achieves a win rate of 58% against the heuristic bot, approaching a level competitive with human.

With external knowledge, we observe that the agent starts to use very special moves at proper time. As an example

shown in **Figure** 8, a steel-type _Klefki_ is vulnerable to the ground-type attack of the opposing _Rhydon_, a ground-type Pokémon. Usually in such a disadvantage, the agent will choose to switch to another Pokémon, however, it chooses to use the move "Magnet Rise", which levitates the user to make it immune to ground-type moves for five turns. As a result, the ground-type attack "Earthquake" of the opposing _Rhydon_ becomes invalid.

### 6.3. Consistent Action Generation

Existing studies (Wei et al., 2022; Yao et al., 2022; Shinn et al., 2023; Bommasani et al., 2021; Hu et al., 2023) show that reasoning and prompting can improve the ability of LLMs on solving complex tasks. Instead of generating a one-shot action, we evaluate existing prompting approaches including Chain-of-Thought (Wei et al., 2022) (CoT), Self-Consistency (Wang et al., 2022) (SC) and Tree-of-Thought (Yao et al., 2023) (ToT). For CoT, the agent initially generates a thought that analyzes the current battle situation and outputs an action conditioned on the thought. For SC (k=3), the agent generates three times of actions and select the most voted answer as the output. For ToT (k=3), the agent generates three action options and picks out the best one evaluated by itself.

_Table 5._ Performance of prompting approaches in battles against the bot.

| Player | Win rate ↑ | Score ↑ | Turn # | Battle # |
|---|---|---|---|---|
| Human | 59.84% | 6.75 | 18.74 | 254 |
| Origin | 58.00% | 6.53 | 18.84 | 100 |
| CoT | 54.00% | 5.78 | 19.60 | 100 |
| SC (k=3) | **64.00**% | 6.63 | 18.86 | 100 |
| ToT (k=3) | 60.00% | 6.42 | 20.24 | 100 |

**Table** 5 presents the comparison results of the original IO prompt generation and three algorithms. Notably, CoT results in a performance degradation by a 6 AP drop in the win rate. In comparison, SC brings a performance improvement, with the win rate surpassing human players. Beyond the results, our greater interest lies in understanding the underlying reasons for these observations.

As introduced in Section 3.2, for each turn there is single action can be taken, which means if the agent chooses to switch yet the opponent choose to attack, the switch-in Pokémon will sustain the damage. Usually switching happens when the agent decides to leverage the type advantage of an off-the-battle Pokémon, and thus the damage taken is sustainable since the switch-in Pokémon is typically type-resistant to the opposing Pokémon's moves. However, when the agent with CoT reasoning faces a powerful opposing Pokémon, its actions become inconsistent by switching to different Pokémon in consecutive turns, which we call _panic switching_. _Panic switching_ wastes chances of taking moves and leading to the defeat. An illustrative example is shown

*Figure 9.* When facing a powerful Pokémon, the agent with CoT switches different Pokémon in three consecutive to elude the battle. This gives the opponent three free turns to quadruple its attack stats and quickly defeat the agent's entire team.

in **Figure** 9: starting from turn 8, the agent chooses to continuously switch to different Pokémon in three consecutive turns, giving the opposing Pokémon three free turns to boost its attack stats to four times and take down the agent's entire team quickly.

*Table 6.* Statistic analysis of panic switching

| Player | Win rate ↑ | Switch rate | CS1 rate | CS2 rate |
|--------|-----------|-------------|----------|----------|
| Origin | 58.00% | 17.05% | 6.21% | 22.98% |
| CoT | 54.00% | 26.15% | 10.77% | 34.23% |
| SC (k=3) | 64.00% | 16.00% | 1.99% | 19.86% |
| ToT (k=3) | 60.00% | 19.70% | 5.88% | 23.08% |

**Table** 6 provides statistical evidence, where CS1 represents the ratio of active switches where the last-turn action is a switch and CS2 rates represent the ratio of active switches here at least one action from the last two turns is a switch, among all active switches, respectively. The higher the CS1 rate, the greater the inconsistency of generation. Obviously, CoT largely increases the continuous switch rate, whereas, SC decreases the continuous switch rate.

Upon examining the thoughts generated by CoT, we observe that the thoughts contain panic feelings: the agent describes how powerful the opposing Pokémon is and the weaknesses of the current Pokémon, and ultimately decides to switch to another Pokémon, as in "*Drapion* has boosted its attack to two times, posing a significant threat that could potentially knock out *Doublade* with a single hit. Since *Doublade* is slower and likely to be knocked out, I need to switch to *Entei* because...". Action generation conditioned on panic thoughts leads the agent to continuously switch Pokémon instead of attacking. In comparison, consistent action generation with SC decreases the continuous switch ratio by independently generating actions multiple times and voting out the most consistent action as shown in **Figure** 5, leading to a higher win rate. The observation is reflecting: when humans face stressful situations, overthinking and exaggerating difficulties lead to panic feelings and paralyze their ability to take actions, leading to even worse situations.
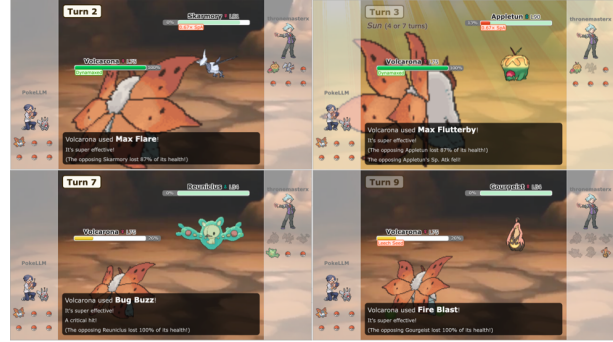


*Figure 10.* POKéLLMON selects effective moves in every turn, causing the opponent's entire team to faint using one Pokémon.

## 7. Online Battle

To test the battle ability of POKéLLMON against human, we set up the eighth-gen battles on Pokémon Showdown, where the agent battled against random human players for the Ladder competitions from Jan. 25 to Jan. 26, 2024. Besides, we invited an human player who has over 15 years of experience with Pokémon games, representing the average ability of human players to play against POKéLLMON.

### 7.1. Battle Against Human Players

*Table 7.* Performance of POKéLLMON against human players.

| v.s. Player | Win rate ↑ | Score ↑ | Turn # | Battle # |
|-------------|-----------|---------|--------|----------|
| Ladder Player | 48.57% | 5.76 | 18.68 | 105 |
| Invited Player | 56.00% | 6.52 | 22.42 | 50 |

**Table** 7 presents the performance of the agent against human players. POKéLLMON demonstrates comparable performance to disciplined Ladder players who have extensive battle experience, and achieves a higher win rate than the invited player. The average number of turns in the Ladder competitions is lower because human players sometimes forfeit when they believe they will lose to save time.

### 7.2. Battle Skill Analysis

**Strength:** POKéLLMON seldom make mistakes at choosing the effective move and switching to another suitable Pokémon due to the KAG strategy. As shown in **Figure** 10, in one battle, the agent uses only one Pokémon to cause the entire opponent team fainted by choosing different attack moves toward different Pokémon.

Moreover, POKéLLMON exhibits human-like attrition strategy: With some Pokémon have the "Toxic" move that can inflict additional damage every turn and the "Recover" move that can recover its HP, the agent starts to first poisoned the opposing Pokémon and frequently uses the "Recover" to prevent itself from fainting. By prolonging the battle, the opposing Pokémon's HP is gradually depleted by the poisoning damage. Using attrition strategy requires an un-
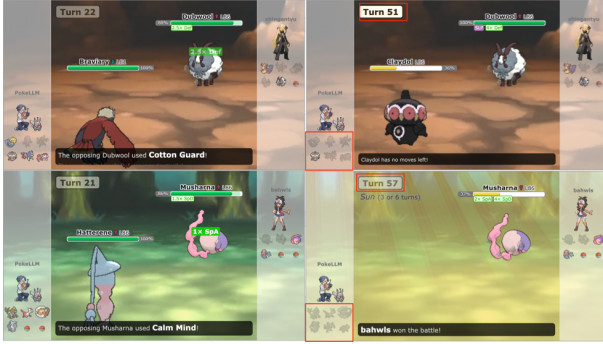
*Figure 11.* POKÉLLMON suffers from attrition strategies: the opponent players frequently recover high-defense Pokémons. Breaking the dilemma requires joint effects across many turns.

derstanding of moves like "Toxic", "Recover" and "Protect", as well as the right timing for their use (such as when there's no type-weakness or when having high defense). An example with battle animation can be found at: https://poke-llm-on.github.io.

**Weakness:** POKÉLLMON tends to take actions that can achieve short-term benefits, therefore, making it vulnerable to human players' attrition strategy that requires long-term effort to break. As shown in the two battles in **Figure** 11, after many turns, the agent's entire team is defeated by the human players' Pokémon, which have significantly boosted defense and engage in frequent recovery. **Table** 8 reports the performance of POKÉLLMON in battles where human players either use the attrition strategy or not. Obviously, in battles without the attrition strategy, it outperforms Ladder players, while losing the majority of battles when human play the attrition strategy.

*Table 8.* Battle performance impacted by the attrition strategy

| Ladder | Win rate ↑ | Score ↑ | Turn # | Battle # |
|---|---|---|---|---|
| w. Attrition | 18.75% | 4.29 | 33.88 | 16 |
| w/o Attrition | 53.93% | 6.02 | 15.95 | 89 |

The "Recover" move recovers 50% HP in one turn, which means if an attack cannot cause the opposing Pokémon more than 50% HP damage in one turn, it will never faint. The key to breaking the dilemma is to firstly boost a Pokémon's attack to a very high stage and then attack to cause unrecoverable damage, which is a long-term goal that requires joint efforts across many turns. POKÉLLMON is weak to the long-term planing because current design does not keep a long-term plan in mind across many timesteps, which will be included in the future work.

Finally, we observe that experienced human players can misdirect the agent to bad actions. As shown in **Figure** 12, our *Zygarde* has one chance to use an enhanced attack move. At the end of turn 2, the opposing *Mawile* is fainted, leading to a forced switch and the opponent choose to switch in



*Figure 12.* An experienced human player misdirects the agent to use a dragon-type attack by firstly sending out a dragon-type Pokémon and immediately switch to another Pokémon immune to the dragon-type attack.

*Kyurem.* This switch is a trick that lures the agent uses a dragon-type move in turn 3 because *Kyurem* is vulnerable to dragon-type attacks. In turn 3, the opponent switches in *Tapu Bulu* at the beginning, a Pokémon immune to dragon-type attacks, making our enhanced attack chance wasted. The agent is fooled because it makes decision only based on the current state information, while experienced players condition on not only the state information, but also the opponent's next action prediction.

Seeing through tricks and predicting the opponent's next action require the agent being disciplined in the real battle environment, which is the future step in our work.

## 8. Conclusion

In this paper, we enable LLMs to autonomously play the well-known Pokémon battles against human. We introduce POKÉLLMON, the first LLM-based agent that achieves human-competent performance in tactical battle games. We introduce three key strategies in the design of POKÉLLMON: (i) In-Context Reinforcement Learning, which consumes the text-based feedback as "reward" to *iteratively* refine the action generation policy without training; (ii) Knowledge-Augmented Generation that retrieves external knowledge to combat hallucination and ensures the agent act timely and properly; (iii) Consistent Action Generation that prevents the *panic switching* issue when encountering powerful opponents. The architecture of POKÉLLMON is *general* and can be adapted for the design of LLM-based agents in many other games, addressing the problems of hallucination and action inconsistency.

Online battles show that POKÉLLMON demonstrates human-like battle ability and strategies, achieving 49% of win rate in the Ladder competitions and 56% of win rate in the invited battles. Furthermore, we uncover its vulnerabilities to human players' attrition strategies and deception tricks, which are considered as our future work.

# References

List of abilities, 2024a. URL https://bulbapedia.bulbagarden.net/wiki/Ability#List_of_Abilities.

List of moves, 2024b. URL https://bulbapedia.bulbagarden.net/wiki/List_of_moves.

List of pokémon by national pokédex number, 2024c. URL https://bulbapedia.bulbagarden.net/wiki/List_of_Pokmon_by_National_Pokdex_number.

Pokémon showdown, 2024. URL https://play.pokemonshowdown.com.

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

Akata, E., Schulz, L., Coda-Forno, J., Oh, S. J., Bethge, M., and Schulz, E. Playing repeated games with large language models. *arXiv preprint arXiv:2305.16867*, 2023.

Bakhtin, A., Brown, N., Dinan, E., Farina, G., Flaherty, C., Fried, D., Goff, A., Gray, J., Hu, H., et al. Human-level play in the game of diplomacy by combining language models with strategic reasoning. *Science*, 378(6624):1067–1074, 2022.

Batra, D., Chang, A. X., Chernova, S., Davison, A. J., Deng, J., Koltun, V., Levine, S., Malik, J., Mordatch, I., Mottaghi, R., et al. Rearrangement: A challenge for embodied ai. *arXiv preprint arXiv:2011.01975*, 2020.

Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

Cabello, L., Li, J., and Chalkidis, I. Pokemonchat: Auditing chatgpt for pok\'emon universe knowledge. *arXiv preprint arXiv:2306.03024*, 2023.

Duan, J., Yu, S., Tan, H. L., Zhu, H., and Tan, C. A survey of embodied ai: From simulators to research tasks. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 6(2):230–244, 2022.

Goertzel, B. Artificial general intelligence: concept, state of the art, and future prospects. *Journal of Artificial General Intelligence*, 5(1):1, 2014.

Goertzel, B. and Pennachin, C. *Artificial general intelligence*, volume 2. Springer, 2007.

Guu, K., Lee, K., Tung, Z., Pasupat, P., and Chang, M. Retrieval augmented language model pre-training. In *International conference on machine learning*, pp. 3929–3938. PMLR, 2020.

Hafner, D., Pasukonis, J., Ba, J., and Lillicrap, T. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.

Hong, S., Zheng, X., Chen, J., Cheng, Y., Wang, J., Zhang, C., Wang, Z., Yau, S. K. S., Lin, Z., Zhou, L., et al. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 2023.

Hu, S., Huang, T., İlhan, F., Tekin, S. F., and Liu, L. Large language model-powered smart contract vulnerability detection: New perspectives. *arXiv preprint arXiv:2310.01152*, 2023.

Hua, W., Fan, L., Li, L., Mei, K., Ji, J., Ge, Y., Hemphill, L., and Zhang, Y. War and peace (waragent): Large language model-based multi-agent simulation of world wars. *arXiv preprint arXiv:2311.17227*, 2023.

Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.

Light, J., Cai, M., Shen, S., and Hu, Z. From text to tactic: Evaluating llms playing the game of avalon. *arXiv preprint arXiv:2310.05036*, 2023.

Ma, W., Mi, Q., Yan, X., Wu, Y., Lin, R., Zhang, H., and Wang, J. Large language models play starcraft ii: Benchmarks and a chain of summarization approach. *arXiv preprint arXiv:2312.11865*, 2023.

Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pp. 1928–1937. PMLR, 2016.

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

Park, J. S., O'Brien, J., Cai, C. J., Morris, M. R., Liang, P., and Bernstein, M. S. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pp. 1–22, 2023.

Patil, S. G., Zhang, T., Wang, X., and Gonzalez, J. E. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*, 2023.

Rawte, V., Sheth, A., and Das, A. A survey of hallucination in large foundation models. *arXiv preprint arXiv:2309.05922*, 2023.

Sahovic, H. Poke-env: pokemon ai in python, 2023a. URL https://github.com/hsahovic/poke-env.

Sahovic, H. poke-env: Heuristicbot, 2023b. URL https://github.com/hsahovic/poke-env/blob/master/src/poke_env/player/baselines.py.

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

Shinn, N., Labash, B., and Gopinath, A. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366*, 2023.

Significant Gravitas. AutoGPT. URL https://github.com/Significant-Gravitas/AutoGPT.

Singh, I., Blukis, V., Mousavian, A., Goyal, A., Xu, D., Tremblay, J., Fox, D., Thomason, J., and Garg, A. Progprompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 11523–11530. IEEE, 2023.

Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

Wang, G., Xie, Y., Jiang, Y., Mandlekar, A., Xiao, C., Zhu, Y., Fan, L., and Anandkumar, A. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023a.

Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., Chen, Z., Tang, J., Chen, X., Lin, Y., et al. A survey on large language model based autonomous agents. *arXiv preprint arXiv:2308.11432*, 2023b.

Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., Chowdhery, A., and Zhou, D. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022.

Wang, Z., Cai, S., Liu, A., Ma, X., and Liang, Y. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. *arXiv preprint arXiv:2302.01560*, 2023c.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35: 24824–24837, 2022.

Wikipedia. Gameplay of pokémon, 2023. URL https://en.wikipedia.org/wiki/Gameplay_of_Pok%C3%A9mon.

Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B., Zhang, M., Wang, J., Jin, S., Zhou, E., et al. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*, 2023.

Xu, Y., Wang, S., Li, P., Luo, F., Wang, X., Liu, W., and Liu, Y. Exploring large language models for communication games: An empirical study on werewolf. *arXiv preprint arXiv:2309.04658*, 2023.

Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., and Cao, Y. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.

Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T. L., Cao, Y., and Narasimhan, K. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*, 2023.