POKÉLLMON: A GROUNDING AND REASONING BENCHMARK FOR LARGE LANGUAGE MODELS IN POKÉMON BATTLES

Anonymous authors

Paper under double-blind review

Abstract

Developing grounding techniques for LLMs poses two requirements for interactive environments, *i.e.*, (i) the presence of rich knowledge beyond the scope of existing LLMs and (ii) the complexity of tasks that require strategic reasoning. Existing environments fail to meet both requirements due to their simplicity or reliance on commonsense knowledge already encoded in LLMs for interaction. In this paper, we present PokéLLMon, a new benchmark enriched with fictional game knowledge and characterized by the intense, dynamic, and adversarial gameplay of Pokémon battles, setting new challenges for the development of grounding and reasoning techniques in interactive environments. Empirical evaluations demonstrate that existing LLMs lack game knowledge and struggle in Pokémon battles. We investigate grounding techniques that leverage game knowledge and self-play experience, and provide a thorough analysis of reasoning methods from a new perspective of action consistency. Additionally, we introduce higher-level reasoning challenges when playing against human players. The implementation of our benchmark is anonymously released at: https://anonymous.4open.science/r/PokeLLMon.

027 028 029

025

026

006

008 009 010

011 012 013

014

015

016

017

018

019

021

1 INTRODUCTION

030

The success of Large Language Models (LLMs) comes from encoding massive textual data by
predicting the next token with huge model capacity, and generalizing well to various tasks (Ouyang
et al., 2022; Brown et al., 2020; Achiam et al., 2023; Xi et al., 2023; Wang et al., 2023b). The
pre-trained LLMs, encode the knowledge from text and exhibit cognitive abilities such as reasoning
and planning, ways of organizing knowledge described in text inherently.

In other words, LLMs lack experiential grounding (Mahowald et al., 2024; Hu et al., 2024), which
prevents them from understanding new concepts outside their scope, or evolving their reasoning
abilities in the interactive environments. Recent research focuses on grounding LLMs in games via
Reinforcement Learning (RL) (Carta et al., 2023; Tan et al., 2024) or supervised fine-tuning (Zhu
et al., 2023; Feng et al., 2023). However, it has been observed that even a well-trained agent still lacks
generalizability to slightly altered settings, such as substituting action tokens with synonyms (Carta
et al., 2023), suggesting that simplistic environments with limited scenarios do not best facilitate the
development of grounding and reasoning techniques.

Developing grounding and reasoning methods for LLMs places two requirements for interactive environments: (i) the presence of knowledge beyond the scope of existing LLMs and (ii) the complexity of tasks that demand strategic reasoning abilities. The environments used in previous work do not meet both of these requirements due to their reliance on commonsense knowledge encoded in LLMs (Shridhar et al., 2020; Xiang et al., 2024) (a locked door can be unlocked with a key), or their simplicity and limited number of scenarios (Carta et al., 2023; Tan et al., 2024).

Scope and Contributions: In this paper, we introduce PokéLLMon, a new grounding and reasoning
 benchmark for LLMs in Pokémon Battles, which sets a new challenge for LLMs to master fictional
 game knowledge that falls outside of their current scope (Cabello et al., 2023). To date, there are over
 1,000 Pokémon species and 900 battle moves (bul, 2024b;a), offering a wealth of game knowledge
 and a large amount of combination possibilities, making the game highly dynamic. Furthermore,

0							
054	Environment	Imperfect Info.	Knowledge	Strategic	Adversarial	Dynamic	Game⇒Text
055 056 057	ALFWorld (Shridhar et al., 2020) ScienceWorld (Wang et al., 2022a)	√ √ ×	Low Low+	Low Low+	× × ×	Low Low	Lossless Lossless
058	OverCooked-AICarroll et al. (2019) Crafter (Hafner, 2021)	×	Low Low Medium	Low+ Medium	××	Low Medium High	Lossiess Lossy Lossy
059	Minecraft (Mojang Studios) StarCraft II (Vinyals et al., 2017)	√ √	High High	Medium High	× √	High High	Lossy Lossy
061	PokéLLMon (Ours)	\checkmark	High	High	\checkmark	High	Lossless

063Table 1: Comparison with popular environments in LLM research across several aspects: Imperfect064Info: The important game information is partially observable. Knowledge: The level of knowledge065that beyond the scope of LLMs; Strategic: The strategic level of playing the game; Human: Whether066human players can be involved; Adversarial: Whether it is an adversarial game or not; Game \Rightarrow Text:067whether translating the game into text is lossy or lossless.

the adversarial feature creates an intense gameplay experience with a high ceiling for reasoning,
 especially in the presence of powerful opponents such as human experts.

This paper conducts comprehensive evaluation of existing LLMs on game knowledge prediction and 071 Pokémon battles. Empirical evaluations demonstrate that LLMs, especially open-source LLMs, suffer 072 from a severe lack of game knowledge and thus struggle to generate effective actions, indicating our 073 benchmark is a good testbed for grounding techniques. To ground LLMs in games, we first evaluate 074 the impact of game knowledge for different LLMs. Experiments show that the improvement of 075 knowledge is highly related to the inherent reasoning abilities of LLMs, suggesting that when the game 076 knowledge is not the bottleneck, reaching a higher-level performance requires good reasoning ability. 077 Further, we evaluate grounding with self-evolution techniques, which shows that the adversarial setting of PokéLLMon offers an ideal curriculum for learning from self-play experience.

In intense adversarial games, action consistency is an important indicator of performance. Through our analysis, we discover that existing reasoning approaches such as Chain-of-Thoughts (Wei et al., 2022) (CoT) and Reflexion (Shinn et al., 2023) can both lead to inconsistent actions, *i.e.*, the LLM frequently switches Pokémon in consecutive steps and wastes chances to attack. To this end, we introduce an approach that recursively refines the thoughts from the previous steps and can significantly enhance action consistency and gameplay performance. Finally, to demonstrate higher-level reasoning challenges, we test the best-performing LLM in battles against invited human players. The evaluation shows that the LLM player struggles to overcome human players' attrition and misdirection strategies, suggesting a significant room for improvement in reasoning approaches.

088

091 092

094

096

In summary, this paper makes three contributions:

- We introduce PokeLLMon, a benchmark that poses new challenges for grounding and reasoning in an environment which features abundant fictional game knowledge and strategic gameplay.
- We conduct a comprehensive evaluation of LLMs' gameplay performance, suggesting that their lack of game knowledge and struggle in Pokémon battles. We further investigate grounding techniques that leverage game knowledge and self-play experience.
- We provide a thorough analysis of reasoning approaches from a new perspective of action consistency, and introduce an effective approach to improve consistency. Additionally, we present high-level reasoning challenges when competing against disciplined human players.
- 097 098 099

100

2 RELATED WORK

Interactive environments: The complexity of environments is depend on two key dimensions:
 the volume of game knowledge and the strategic depth required. Table 1 summarizes popular
 environments *w.r.t.* several aspects. Environments that are widely used in recent LLM research were
 introduced before the rise of LLMs, such as ALFWorld (Shridhar et al., 2020), ScienceWorld (Wang
 et al., 2022a), and BabyAI (Maxime et al., 2018), featuring limited scenarios or fixed task settings
 and mainly relying on commonsense knowledge to play. Although the open-world characteristics of
 Crafter (Hafner, 2021) and Minecraft (Mojang Studios) make them good sandboxes for exploration
 and crafting tasks, they are not designed as intense adversarial games that require strategic reasoning.



125 Figure 1: Illustration of how LLMs interact with the environment: At the current time step t, the 126 environment outputs an observation prompt o_t that describes the observable information of the 127 current battle state, and the feedback f_{t-1} that describes the action execution outcome of the last 128 time step. The LLM then takes o_t and previous f_i (i < t) as input, conducts reasoning to formulate strategies, and returns an action a_t to the environment. Within the environment, a client is response 129 for translating the text-based o_t and a_t into symbolic information and communicate with the game 130 engine, while the game engine is responsible for generate the next game state based on the current 131 game state and received actions from both players. 132

StarCraft II (Vinyals et al., 2017), a real-time strategy game, presents challenges for LLMs due to the demand for intensive controls and the difficulty of representing vision-based game states in text (Ma et al., 2023). In comparison, PokéLLMon, which encompasses a high volume of game knowledge and offers strategic gameplay in a format friendly to LLMs (does not require intense control and can be directly translated into text), is an ideal testbed for LLM research.

LLMs in games: There are primarily three categories of LLM game agents: (1) Prompt-based 139 approaches that leverage the reasoning abilities of LLMs and feedback from environments, enabling 140 LLMs to iteratively refine strategies. ReAct (Yao et al., 2022) (ALFWorld) introduces the thinking 141 step as a proxy for sub-tasks. Reflexion (Shinn et al., 2023) (ALFWorld) and DEPS (Wang et al., 142 2023d) (Minecraft) generate self-reflection/explanation based on failure signals and reuse these 143 thoughts for the next trial; In Minecraft, Voyager (Wang et al., 2023a), JARVIS-1 (Wang et al., 2023c) 144 and GTIM (Zhu et al., 2023) iteratively re-generate action code using error messages; (2) Supervised 145 fine-tuning-based methods that collect high-quality trajectories to fine-tune LLMs: E2WM (Xiang 146 et al., 2024) collects embodied experience (VirtualHome (Puig et al., 2018)) using Monte Carlo Tree 147 Search; LLAMARider (Feng et al., 2023) (Minecraft) gathers experience through self-reflection with feedback; (3) Reinforcement learning: GLAM (Carta et al., 2023) (BabyAI-Text) and TWO-148 SOME (Tan et al., 2024) (OverCooked-AI) discipline LLMs using the PPO algorithm (Schulman 149 et al., 2017). 150

151 152

153

133

3 Environment

As shown in Figure 1, the environment of PokéLLMon follows a client-server architecture, where the client is implemented as a text interface for LLMs to perceive the game, and the server is an open-source game engine ¹ for game execution. The game is synchronized at every step: In each step, both players receive observations and select actions synchronously, and the game engine processes these actions to calculate the next step of the game state. In a battle, each player sends out one Pokémon onto the field, keeping the others off the field for potential switches. The winning condition is to make all the opponent's Pokémon faint by reducing their Hit Points (HP) to zero.

¹https://github.com/smogon/Pokémon-showdown, licensed under the MIT License.

Choice	A. Super ef	fective (51	samples)	C. Ineffe	ctive (59 sat	mples)	D. No e	ffect (8 sam	ples)	Overall
LLM	Precision	Recall	\mathbf{F}_1	Precision	Recall	\mathbf{F}_1	Precision	Recall	\mathbf{F}_1	\mathbf{F}_1
Mistral-7B	0.5000	0.0588	0.1053	0.2190	0.3770	0.2771	0.0000	0.0000	0.0000	0.1841
Gemma-7B	0.1672	1.0000	0.2865	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.1238
LLaMA3-8B	0.1818	0.1569	0.1684	0.1871	0.5246	0.2759	0.1333	0.2500	0.1739	0.2225
LLaMA2-7B	0.1589	1.0000	0.2742	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.1185
LLaMA2-13B	0.1721	0.4118	0.2428	0.1500	0.3443	0.2090	0.0000	0.0000	0.0000	0.2094
LLaMA2-70B	0.1902	0.7647	0.3047	0.1613	0.1639	0.1626	0.0000	0.0000	0.0000	0.2130
GPT-3.5-turbo	0.3778	0.3333	0.3542	0.1944	0.8033	0.3131	0.3333	0.2500	0.2857	0.3290
GPT-4	0.9787	0.9020	0.9388	0.7273	0.7869	0.7559	0.7273	1.0000	0.8421	0.8408
GPT-40	0.9804	0.9804	0.9804	0.7000	0.8033	0.7481	0.7273	1.0000	0.8421	0.8549
GPT-4o-mini	0.5652	0.2549	0.3514	0.3750	0.5902	0.4586	0.4000	0.7500	0.5217	0.4172

Table 2: Evaluation of LLMs on type effectiveness prediction

173

162

174 **Observation:** Pokémon battles are an imperfect information game, suggesting that the game state 175 is partially observable. The observation includes the information (stats, abilities, and moves) of a 176 player's own Pokémon team and partial information of the opponent's team, *i.e.*, the species and HP stats of appeared Pokémon. Although the observation is represented as an image from human 177 perspective, it can be symbolically represented without loss of information. At time step t, the client 178 translates a symbolic observation into a text observation o_t , as shown in Figure 1. 179

Action: There are two types of actions can be chosen: (1) use one of four battle moves, where a move 181 can cause instant damage or produce special effects. The priority of taking moves is determined by the current speed of two Pokémon on the field. (2) Switch to an off-the-field Pokémon. If choose to 182 switch, the switching will happen immediately before the opposing Pokémon take a move, and the 183 switch-in Pokémon cannot take an extra action in this step. The admissible actions (up to 9 choices, 184 *i.e.*, 4 moves plus 5 switch-in options) are included in o_t . The LLM is adopted as the policy model π_{θ} 185 to generate reasoning and an action a_t .

187 Feedback: The action will then be executed and during execution, human players perceive the feedback from the battle animation that reflects the effectiveness of chosen actions. To compensate 188 this information, we introduce four types of textual feedback: (1) The change in HP. (2) The outcome 189 of actions in terms of type-effectiveness, *i.e.*, whether it is super-effective, ineffective, or has no effect. 190 (3) The priority of move execution. (4) The effects of special moves on stats/status changes, weather, 191 and side conditions. The feedback f_i (i < t) is included in the o_t for LLMs to perceive its previous 192 actions and outcomes. In the environment, we also provide numeric reward value derived based on 193 the outcome of actions. 194

4 PRELIMINARY EVALUATION

4.1 GAME KNOWLEDGE EXAMINATION

Type-effectiveness is fundamental knowledge in Pokémon battles. It defines the effectiveness of a certain type of Pokémon when attacked by a certain type of attack. For example, a Fire-type Pokémon is vulnerable to Water-type attacks (see the chart in Appendix D). We use the chart to generate multiple-choice questions to test the game knowledge of LLMs. The question template is as follows:

Multi-choice question: In Pokémon battles, a type₁ attack is against a type₂ Pokémon. A. Super-effective (2x) B. Standard (1x) C. Ineffective (0.5x) D. No effect (0x)

206 Table 2 shows the results of existing LLMs, where we report the precision, recall, and F1 score for 207 choices A, C, and D (with B being the majority answer), as well as the weighted F1 score across 208 the three choices. We observe that the overall F1 score is low, especially for open-source LLMs, 209 suggesting that game knowledge is largely beyond the scope of pre-trained LLMs. Even though 210 the best-performing LLM, GPT-40, achieves quite accurate results, we will show that this is still 211 insufficient in intense battles where a single mistake can lead to a significant disadvantage.

212 213

214

195

196

197

199

200

201

202 203

204

205

4.2 GAME PERFORMANCE EVALUATION

Table 3 shows the gameplay performance of existing LLMs against ExperSystem, an expert system 215 that simulates human's decision-making process with numeric damage calculation (details in Ap-

Shot #		<i>n</i> =0			n=1			<i>n</i> =3	
Player	Win Rate↑	Battle Score↑	Error Rate↓	Win Rate↑	Battle Score↑	Error Rate↓	Win Rate↑	Battle Score↑	Error Rate↓
ExpertSystem	0.5000	6.000	0.000	-	-	-	-	-	-
Human	0.5984	6.750	0.000	-	-	-	-	-	-
Random	0.0120	2.340	0.000	-	-	-	-	-	-
MaxPower	0.1040	3.790	0.000	-	-	-	-	-	-
Mistral-7B	0.0486	3.047	0.008	-	-	>0.90	-	-	>0.90
Gemma-7B	0.0866	3.684	0.022	-	-	>0.90	-	-	>0.90
LLaMA3-8B	0.1034	3.902	0.028	-	-	>0.90	-	-	>0.90
LLaMA2-7B	0.0760	3.572	0.015	-	-	>0.90	-	-	>0.90
LLaMA2-13B	0.0829	3.610	0.006	-	-	>0.90	-	-	>0.90
LLaMA2-70B	0.1065	3.703	0.031	-	-	>0.90	-	-	>0.90
GPT-3.5-turbo	0.1351	4.089	0.000	0.0988	3.673	0.000	0.1333	4.027	0.000
GPT-40-mini	0.0719	3.127	0.000	0.1236	3.789	0.000	0.1634	4.254	0.000
GPT-40	0.3100	4.719	0.000	0.2917	4.605	0.000	0.2986	4.578	0.001
GPT-4	0.2673	4.954	0.000	0.2644	4.926	0.000	0.2692	4.603	0.000

Table 3: Evaluation of the gameplay performance of LLMs (*n* is the number of few-shot examples)

pendix A). In each step, LLMs take o_t as input and output a_t without explicit reasoning (see Section 6 for the evaluation of reasoning). n is the number of few-shot examples used in prompt. To provide more comparisons, we report the performance of Human (randomly-matched human players from game server²), Random (method that randomly selects actions) and Max-Power (method that selects moves with the highest move power). Each experiment is run over 1,000 times for open-source LLMs, and over 200 times for GPTs. The error rate represents the proportion of inadmissible actions generated by LLMs. The battle score is defined as follows:

Battle Score =
$$\sum_{p \in \mathcal{P}} \operatorname{HP}(p) + \sum_{p \in \mathcal{O}} (1 - \operatorname{HP}(p))$$
 (1)

where HP(p) is the percentage of Pokémon p's HP at the end of a battle, \mathcal{P} and \mathcal{O} are the player's team and the opponent's team. The score is the sum of the remaining HP percentages of the player's Pokémon and the HP loss percentages of the opponent's Pokémon, rewarding both preservation of health and infliction of damage. For LLMs, the temperature τ is set to 0 to reduce inconsistency.

244 From Table 3, we make three observations: (1) all the LLMs, especially open-source LLMs, perform 245 badly in Pokémon battles, and their gameplay performance is positively correlated with their perfor-246 mance in type-effectiveness prediction. (2) With few-shot examples, open-source LLMs repeat the 247 actions in examples rather than generation, making the output inadmissible (error rate >90%). We 248 conjecture that is because they are less aligned with human instructions compared to GPTs, making 249 them difficult to generalize to cases they have never seen. (3) Few-shot examples do not bring any 250 improvement for GPT-40 and GPT-4. This is because that Pokémon battles are dynamic and involve 251 numerous possibilities, making it difficult for few-shot examples to cover all situations.

5 GROUNDING

In this section, we introduce two ways of grounding LLMs in games: grounding with knowledge and with self-play experience.

5.1 GROUNDING WITH KNOWLEDGE

In PokéLLMon, we provide interfaces for knowledge retrieval during battles. The game knowledge is crawled from Pokémon Wiki (pok) and Bulbapedia (bul), and structured in key-value pairs. Specifically, we adopt two essential categories of game knowledge: (i) Type-effectiveness, which describes the strengths and weaknesses of both the attack type and the Pokémon type, for example:

Key: Fire-type Pokémon **Value:** Is resistant to Fire, Grass, Ice attacks, and is vulnerable to Water, Ground, Rock attacks.

(ii) Move effect, which describes the effects of moves, for example:

267 268 269

237

238 239

252 253

254 255

256 257

258 259

260

261

262

263 264

265

²Playing with humans follows the bot usage policy of the game content provider.

*

270 271 272

273

274

275

276

277

(1) LLM uses the Toxic move

to poison the opposing Pokemon

278

279

280 281 282

283

284

285

287

288

289

295

Key: Toxic Value: Toxic poisons the target, causing it to lose progressively increasing HP each turn.

(2) The LLM frequently uses the Recover

move to prevent its Pokémon from fainting

Figure 2: With game knowledge, the LLM exhibits an attrition strategy using two moves: It first

uses *Toxic* to poison the opponent, which inflicts additional poisoning damage on every turn. Then, it

prolongs the battle by frequently healing itself with *Recover*, a move that can restores 50% of HP. As a result, the opponent gradually weakens due to the poisoning damage and faints after 7 turns.

Turn 23

(3) The opposing Pokémon is

depleted by the poisoning damage

*

During games, the agent retrieves the type-effectiveness knowledge using the opponent Pokémon's type, and move knowledge using the move name of its own Pokémon. Retrieved knowledge will be integrated into the observation o_t as the input for LLM to generate action a_t .

Tuble 1. Evaluation on the impact of game knowled	50
Table 4. Evaluation on the impact of same knowled	oe.

LLM	Win Rate	W.R. w/ know.	Battle Score	B.S. w/ know.	\mathbf{F}_1
GPT-3.5-turbo	0.1351	0.0795 (-5.56%)	4.089	3.927 (-0.162)	0.3290
GPT-4o-mini	0.0719	0.1584 (+8.65%)	3.127	3.695 (+0.568)	0.4172
GPT-40	0.3100	0.4217 (+11.17%)	4.719	5.413 (+0.694)	0.8549
GPT-4	0.2673	0.4744 (+20.71%)	4.954	5.869 (+0.915)	0.8408

296 In Table 4 we evaluate the impact of game knowledge in battles against ExpertSystem, where F_1 297 is the overall F1 score of type-effectiveness prediction task. Among four LLMs, game knowledge 298 significantly enhances GPT-4, with an increase of 20.71% in win rate, while it even hurts the 299 performance of GPT-3.5-turbo, leading to a win rate drop of 5.56%. This suggests that improvement 300 in knowledge is related to the inherent reasoning ability of LLMs. Moreover, by comparing GPT-40 301 with GPT-4, we can conclude that when game knowledge is not the bottleneck, achieving higher-level 302 performance requires good (explicit/implicit) reasoning ability.

303 With knowledge, we observe that GPT-4 exhibits emergent behaviors: As shown in Figure 2, when 304 the agent has a move that can inflict additional damage regularly, such as Toxic, and another move 305 that can recover its HP, such as *Recover*, the LLM develops the attrition strategy: It first poisons the 306 opposing Pokémon to cause regular damage every turn and then frequently recovers HP to prevent 307 its Pokémon from fainting. After 7 turns, the opponent's HP is gradually depleted by the poisoning 308 damage until it faints.

309 310

5.2 GROUNDING WITH SELF-PLAY EXPERIENCE

311 Adversarial games offer an ideal learning curriculum by allowing self-play against an opponent of the 312 same gameplay level, as demonstrated in AlphaGoZero (Silver et al., 2017) and OpenAI Five (Berner 313 et al., 2019b). We introduce grounding LLMs through learning from self-play experience. Self-314 evolution can be considered as iterations of two phases: trajectory sampling and fine-tuning. For 315 trajectory sampling, we enable two agents initialized from the same LLM to act as a pair of opponents 316 in Pokémon battles to collect rollouts. In the fine-tuning stage, we adopt two approaches that can be 317 easily integrated into self-play: 318

- 319 • Rejection sampling Fine-Tuning (RFT) (Yuan et al., 2023): RFT samples the trajectories of winners to fine-tune LLMs by maximizing the log-likelihood of actions sampled from winning 320 trajectories. The loss function can be formalized as:
- 321 322

 $L = -\frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \log \pi_{\theta}(a_{i,t}^{w} | o_{i,t}^{w})$ (2)

Tag	Method	Win Rate \uparrow	Battle Score \uparrow
Ο	Origin	0.1075	3.904
	RFT	0.1161	4.194
	DPO	0.1212	4.207
C	Origin	0.5641	6.003
	RFT	0.5920	6.266
	DPO	0.5984	6.302

Table 5: Evaluation of self-evolution techniques. ^(D) denotes playing against ExpertSystem and ^(C) denotes playing against MaxPower.

where $(o_{i,t}^w, a_{i,t}^w) \in \tau_i^w$, and τ_i^w is the trajectory of the winner in *i*-th battle.

• Direct Preference Optimization (DPO) (Rafailov et al., 2024): For the *i*-th battle, we sample pairs of $(o_{i,t}^w, a_{i,t}^w)$ and $(o_{i,t}^l, a_{i,t}^l)$ from winner trajectory τ_i^w and loser trajectory τ_i^l , then fine-tune π_{θ} by minimizing the following loss function:

$$L = -\frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(a_{i,t}^{w} | o_{i,t}^{w})}{\pi_{\text{ref}}(a_{i,t}^{w} | o_{i,t}^{w})} - \beta \log \frac{\pi_{\theta}(a_{i,t}^{l} | o_{i,t}^{l})}{\pi_{\text{ref}}(a_{i,t}^{l} | o_{i,t}^{l})} \right) \right]$$
(3)

where π_{ref} is the reference model initialized using the π_{θ} from the last iteration to prevent excessive drift. Optimizing the DPO loss is essential to increase the likelihood margin between the winner actions and loser actions.

345 The implementation details can be found at Appendix A. From Table 5 we observe that both RFT 346 and DPO obtains better performance than the original LLM player, suggesting that LLMs are able to 347 learn in environments with self-play experience. The second observation is that DPO outperforms 348 RFT. This is because RFT only learns from the trajectories of winners, while DPO minimizes the 349 likelihood of losers' actions, making the learning more effective. Overall, adversarial self-play makes self-evolution efficient: If it is not a self-play setting, RFT will be difficult to obtain winner 350 trajectories when playing against a powerful opponent, or quickly get saturated when playing against 351 a very weak opponent. 352

353 354 6 REASONING

324

327 328

334

335

336

342

343

344

356

368

369 370

371 372

6.1 EVALUATION OF REASONING METHODS

In intense adversarial games, inconsistent actions often lead to defeat by wasting opportunities to make strategic moves. To benchmark existing reasoning approaches, we introduce two additional metrics that measure the degree of action inconsistency and are closely linked to gameplay performance.

Metrics: In Pokémon Battles, there are two categories of steps: (1) force switch step that is compulsory to switch when the pokemon on-the-field faints and the player decides which pokemon off-the-field to switch in; (2) active step, *i.e.*, to take move or switch Pokemon. Switching in an active step provides advantages, such as having a Pokémon with type effectiveness over the opposing Pokémon, however, it also means to give the opposing Pokemon a free turn to take a move (and abandon the boost of stats if have). Frequent switching can waste opportunities to attack and will lead to defeat. Therefore, we use switch rate and consecutive switch rate to measure the action consistency of a player, formalized as follows:

Switch Rate =
$$\frac{\text{\# of active switch}}{\text{\# of active step}}$$
 (4)

Consecutive Switch Rate =
$$\frac{\text{# of consecutive active switch}}{\text{# of active switch}}$$
 (5)

Switch rate measures the proportion of times a player makes an active switch, and consecutive switch rate measures the frequency with which a player switches in consecutive steps against the same opponent.

9

Approaches: As shown in Figure 3, we evaluate reasoning approaches, including IOPrompt, Chain-of-Thoughts (Wei et al., 2022), Self Consistency (Wang et al., 2022b) (SC-CoT), Tree-of-Thoughts (Yao et al., 2023), Relexion (Shinn et al., 2023) and Last-Thoughts (See Appendix A for more descriptions).

388

431



Figure 3: Illustration of reasoning approaches evaluated in Table 6. Table 6: Evaluation of reasoning approaches *w.r.t.* gameplay performance and action consistency.

Method	Win Rate↑	Battle Score ↑	Switch Rate	Con. Switch Rate↓
IOPrompt	0.4217	5.413	0.3356	0.2442
CoT	0.3713	5.127	0.3344	0.2647
SC-CoT	0.4065	5.037	0.3643	0.0954
ТоТ	0.2549	4.398	0.3163	0.2938
Reflexion	0.2923	5.015	0.3680	0.2982
LastThoughts	0.4667	5.840	0.2227	0.0861
IOPrompt (τ =0.0)	0.4217	5.413	0.3356	0.2442
IOPrompt (τ =0.5)	0.3818	5.285	0.3204	0.2504
IOPrompt (τ =1.0)	0.3019	4.937	0.3224	0.2689

We adopt GPT-40³ as the LLM and game knowledge, and run every experiments over 200 trials. Two one-shot examples (one for active step and one for force switch step) are designed to guide reasoning. Temperature τ is set to 0 to reduce inconsistency besides for SC-CoT, which requires to encourage different reasoning, we set τ to 0.5. For SC-CoT and ToT, the number of branch *b* is set to 3.

Analysis: Table 6 reports the gameplay performance with reasoning approaches. Surprisingly, we observe that approaches like CoT, ToT and Reflexion significantly decreases the win rate and battle score compared to the vanilla IOPrompt without any reasoning. According to the metrics, the drop of game performance can be attributed to the increase in consecutive switch rate, suggesting that with reasoning, the agent is tend to switch in a new pokemon and switch it out in the next turn, without taking moves. Below, we break down the analysis for every approach.

CoT: Why does CoT lead to action inconsistency? Let us consider two cases: In the first case, 410 our Pokémon is facing an opponent, and in the second case, everything is the same except that the 411 opponent has boosted its attack stats twofold. For IOPrompt, the difference between two cases is 412 only two extra tokens indicating the boosted stats. However, with reasoning, the LLM will describe 413 the disadvantage of facing a boosted opponent in its thoughts. From a generation perspective, the 414 thoughts increase the discrepancy between two observations in the representation space, making the 415 LLM more likely to generate different actions. From a gameplay perspective, conditioned on these 416 thoughts, the LLM tends to protect its Pokémon from fainting by switching it out of the battlefield, 417 yet neglects the fact that consecutive switching actually wastes the chance to attack. 418

SC-CoT and ToT: As show in Table 6, Compared to CoT, SC-CoT decreases the consecutive switch 419 rate while increasing the switch rate, which suggests that SC-CoT is more consistent in consecutive 420 steps with similar observations, yet more likely to switch Pokémon when facing different opposing 421 Pokémon. This is because the majority voting reduces the inconsistency brought by the randomness 422 of sampling. However, it enhances the probability of switching when switching is the predominant 423 option in the probabilistic distribution conditioned on thoughts. A competitor, ToT, replacing the 424 majority voting of SC-CoT by self-evaluation, although outperforms SC-CoT in easily-evaluated tasks 425 like Game of 24 (Yao et al., 2023), works even worse than CoT and SC-CoT in our benchmark. This is because the LLM is unable to make correct self-evaluation due to its lack of game experience, and 426 thus tends to agree with the proposals, which does not provide any benefit for reducing inconsistency, 427 and sometimes even prioritizes the sub-optimal actions, leading to a drop of performance. 428

Last-Thoughts: We introduce a simple yet effective solution: Last-Thoughts. By recursively taking
 the thoughts from the last step into current reasoning procedure, we observe a significant drop of

³GPT-40 is more cost-efficient than GPT-4, allowing us to run more experimental trials for accurate estimation.

432 consecutive switch rate from 34.50% to 9.31%, which thereby boosts the win rate from 35.00% to 433 46.67%. With Last-Thoughts, the LLM refines its thoughts based on its previous thoughts instead of 434 generating from scratch. For the perspective of generation, last thoughts can be deemed as a summary 435 of the last observation, thus two consecutive observations become more similar in the representation 436 space, leading to consistent outputs.

437 **Reflexion:** As shown in Figure 3, the LLM uses reflections on the previous step to generate a new 438 action. There are two reasons that reflection does not work well: (1) unlike static task settings (Shrid-439 har et al., 2020), Pokémon battles are dynamic, suggesting that reflections from the past are likely 440 outdated; (2) Even if not outdated, the LLM tends to think that the last action did not meet its 441 expectations, and thus chooses to switch to a new Pokémon to attack, leading to the increase of 442 inconsistency. However, we believe reflection can be beneficial in a dynamic environment if adjusted appropriately, *i.e.*, by generating reflections offline and retrieving them based on similarity during 443 battles. As a result, the problem of outdated reflections can be addressed, and the agent also avoids 444 wasting chances in a trial-and-error style. 445

446 **Impact of temperature on consistency**: Temperature τ controls the sharpness of distribution for 447 sampling tokens. In Table 6 we report the performance of IOPrompt w.r.t. temperatures varying from 448 0 to 2, suggesting that a lower temperature can reduce inconsistency and increase performance.

Findings: In an intense adversarial game, action consistency is an important indicator for gameplay performance. The inconsistency introduced by CoT is due to the increased discrepancy between consecutive observations, and can be reduced by integrating the previous thoughts in the next step reasoning (Last-Thoughts). Furthermore, a lack of game experience and the highly dynamic feature undermine approaches relying on self-evaluation such as ToT and Reflexion.

6.2 REASONING CHALLENGES IN ONLINE BATTLES

Ladder Player

458 Due to the characteristics of adversarial games, the difficulty of reasoning can be further increased 459 by playing against powerful opponents, e.g., human experts. In order to demonstrate higher-level 460 reasoning challenges, we set up online battles with human players randomly matched from the game 461 server. Online games adhere to the bot usage principles of the game content provider and follow the principle of disclosure: human players are informed through the invitations that they are playing 462 against a non-human agent, and battles only begin upon acceptance. 463

464 465

449 450

451

452

453

454 455 456

457

466

467

471

473

Table 7: Performance on online battles against human players. v.s. Player Win Rate \uparrow Battle Score **↑** 0.4857

5.76

468 Table 7 shows the results of an enhanced agent ⁴ over 100 online battles against human players 469 matched from the ladder competition. We then introduce two reasoning challenges observed when 470 playing against human players, namely long-term dilemma and Theory-of-Mind inference.

472 6.2.1 LONG-TERM DILEMMA

We observe that the current LLM agent still lacks of long-term planning ability, *i.e.*, it tends to 474 achieve short-term goals like instant damage, and thus is vulnerable to the dilemma that requires 475 long-term strategy to break. A good case to measure the long-term planning ability is the human 476 players' attrition strategy, which frequently recovers the Pokémon's HP and gradually deplete the 477 opponent's HP with occasional attacks or regular damage moves such as Toxic. 478

Table 8 reports the performance of the agent in battles where human players either use the attrition 479 strategy or not. We can observe that the agent lost the majority of games when humans performed the 480 attrition strategy, where, the turn number of battles is significantly increased. The key to breaking the 481 dilemma is to form a plan that cross multiple steps: firstly boost its Pokémon's attack to a high stage 482 and then attack to cause unrecoverable damage, which is a long-term goal that requires joint efforts 483 across many steps. 484

⁴To make the agent capable of playing against human, we include more knowledge including ability/item/-485 condition, and adopt GPT-4 as the policy, making it reaches a win rate of 60% against ExpertSystem.



Figure 4: An experienced human player misdirects the LLM to waste an attack chance: In a force switch, the opponent switched in a Pokémon that has a type weakness to agent's current Pokémon's Dragon-type attack. Naturally, the agent chose to use a Dragon-type attack, while the opponent chose to switch in another Pokémon that is immune to Dragon-type attacks. Since the switching occurs before the attack, the Dragon-type attack chosen by the agent has no effect on the opponent's switch-in Pokémon, resulting in a waste of an enhanced attack chance.

503 504 6.

498

499

500

501

502

6.2.2 THEORY-OF-MIND INFERENCE

Theory-of-Mind (ToM) (Frith & Frith, 2005) thinking involves inferring others' intentions from a shifted perspective, which demonstrate enhancement in imperfect information games (Guo et al., 2023) and cooperation games (Zhang et al., 2023; Agashe et al., 2023). In battles, we observe experienced human players misdirected the LLM to bad actions, because it lacks of ToM thinking and makes decisions solely based on the observation.

As shown in Figure 4, the Pokémon from our side has one chance to use an enhanced attack move. At the end of turn 2, the opposing *Mawile* faints, leading to a force switch, and the opponent player chooses to switch in another Pokémon for the next turn. This is a trick to lure the agent to use a dragon-type move, given that dragon-type attack is super effective to the opposing Pokémon. As expected, the agent chooses to use the dragon-type move. However, before our Pokémon attacks, the opponent switches in another Pokémon that is immune to dragon-type attacks, causing our enhanced attack opportunity wasted.

To summarize, attrition and misdirection strategies by human players pose higher-level reasoning challenges for LLMs. Due to the adversarial gameplay setting and the presence of powerful opponents, the benchmark presents a high ceiling for strategic reasoning.

7 CONCLUSION

This paper introduces a new grounding and reasoning benchmark for LLMs in Pokémon Battles, 523 featuring rich game knowledge and a high ceiling for strategic reasoning due to its highly dynamic, 524 intense and adversarial gameplay, especially in the presence of professional opponents. Through 525 experiments we verify the lackness of game knowledge and experience of existing LLMs, investigate 526 the potential grounding solutions with game knowledge and self-play experience. With a thorough 527 analysis, we identify the weakness of existing reasoning approaches from a new perspective of action 528 consistency, and provide an effective solution to reduce inconsistency. Finally, we introduce two 529 advanced strategies exhibited by human players that pose higher-level reasoning challenges. Overall, 530 we believe PokéLLMon is an ideal testbed for developing grounding and reasoning techniques.

531 532

533

521

522

8 ETHICS STATEMENT

The goal of the paper is only for AI development akin to previous benchmarks in StarCraft II (Vinyals et al., 2017) and DOTA (Berner et al., 2019a). We acknowledge that our developments follows the bot usage principle of the service provider⁵. To further address ethical concerns, the environment is implemented in compliance with the principle of disclosure, *i.e.*, for online games, human players are informed that they are playing with non-human players, and games only begin upon acceptance.

⁵https://gist.github.com/Kaiepi/becc5d0ecd576f5e7733b57b4e3fa97e

540 REFERENCES

556

569

570

571

572

576

577

580

581

- 542 Bulbapedia: The community-driven pokémon encyclopedia. https://bulbapedia.
 543 bulbagarden.net/wiki/Main_Page.
- 544 Pokemon wiki. https://pokemon.fandom.com/wiki/Pokmon_Wiki.
- 546 List of moves, 2024a. URL https://bulbapedia.bulbagarden.net/wiki/List_of_ 547 moves.
- 548 List of pokémon by national pokédex number, 2024b. URL https://bulbapedia. bulbagarden.net/wiki/List_of_Pokmon_by_National_Pokdex_number. 550
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
 <u>arXiv preprint arXiv:2303.08774</u>, 2023.
- Saaket Agashe, Yue Fan, and Xin Eric Wang. Evaluating multi-agent coordination abilities in large
 language models. arXiv preprint arXiv:2310.03903, 2023.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. <u>arXiv preprint arXiv:1912.06680</u>, 2019a.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Debiak, Christy
 Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale
 deep reinforcement learning. arXiv preprint arXiv:1912.06680, 2019b.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. <u>Advances in neural information processing systems</u>, 33:1877–1901, 2020.
- Laura Cabello, Jiaang Li, and Ilias Chalkidis. Pokemonchat: Auditing chatgpt for pok\'emon universe knowledge. <u>arXiv preprint arXiv:2306.03024</u>, 2023.
 - Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca Dragan. On the utility of learning about humans for human-ai coordination. Advances in neural information processing systems, 32, 2019.
- Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre-Yves
 Oudeyer. Grounding large language models in interactive environments with online reinforcement
 <u>arXiv preprint arXiv:2302.02662</u>, 2023.
 - Yicheng Feng, Yuxuan Wang, Jiazheng Liu, Sipeng Zheng, and Zongqing Lu. Llama rider: Spurring large language models to explore the open world. arXiv preprint arXiv:2310.08922, 2023.
- ⁵⁷⁸ Chris Frith and Uta Frith. Theory of mind. <u>Current biology</u>, 15(17):R644–R645, 2005.
 - Jiaxian Guo, Bo Yang, Paul Yoo, Bill Yuchen Lin, Yusuke Iwasawa, and Yutaka Matsuo. Suspicionagent: Playing imperfect information games with theory of mind aware gpt-4. <u>arXiv preprint</u> arXiv:2309.17277, 2023.
- Danijar Hafner. Benchmarking the spectrum of agent capabilities. <u>arXiv preprint arXiv:2109.06780</u>, 2021.
- Sihao Hu, Tiansheng Huang, Fatih Ilhan, Selim Tekin, Gaowen Liu, Ramana Kompella, and Ling Liu.
 A survey on large language model-based game agents. arXiv preprint arXiv:2404.02039, 2024.
- Weiyu Ma, Qirui Mi, Xue Yan, Yuqiao Wu, Runji Lin, Haifeng Zhang, and Jun Wang. Large language models play starcraft ii: Benchmarks and a chain of summarization approach. <u>arXiv preprint</u> arXiv:2312.11865, 2023.
- Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and
 Evelina Fedorenko. Dissociating language and thought in large language models. <u>Trends in</u>
 Cognitive Sciences, 2024.

594 595 596	Chevalier-Boisvert Maxime, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu Nguyen, and Yoshua Bengio. Babyai: A platform to study the sample efficiency of grounded language learning. <u>arXiv preprint arXiv:1810.08272</u> , 2018.
597 598	Mojang Studios. Minecraft. https://www.minecraft.net/en-us.
599 600 601 602 603	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. <u>Advances in Neural Information Processing Systems</u> , 35: 27730–27744, 2022.
604 605 606	Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba. Virtualhome: Simulating household activities via programs. In <u>Proceedings of the IEEE conference</u> on computer vision and pattern recognition, pp. 8494–8502, 2018.
607 608 609	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. <u>Advances in Neural Information Processing Systems</u> , 36, 2024.
610 611 612	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. <u>arXiv preprint arXiv:1707.06347</u> , 2017.
613 614	Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic memory and self-reflection. <u>arXiv preprint arXiv:2303.11366</u> , 2023.
615 616 617 618	Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. <u>arXiv</u> preprint arXiv:2010.03768, 2020.
619 620 621	David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. <u>nature</u> , 550(7676):354–359, 2017.
623 624 625	Weihao Tan, Wentao Zhang, Shanqi Liu, Longtao Zheng, Xinrun Wang, and Bo An. True knowledge comes from practice: Aligning llms with embodied environments via reinforcement learning. <u>arXiv</u> preprint arXiv:2401.14151, 2024.
626 627 628	Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, et al. Starcraft ii: A new challenge for reinforcement learning. <u>arXiv preprint arXiv:1708.04782</u> , 2017.
629 630 631 632	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. <u>arXiv</u> preprint arXiv:2305.16291, 2023a.
633 634 635	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. <u>arXiv preprint arXiv:2308.11432</u> , 2023b.
636 637 638	Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. Scienceworld: Is your agent smarter than a 5th grader? <u>arXiv preprint arXiv:2203.07540</u> , 2022a.
639 640 641	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh- ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. <u>arXiv preprint arXiv:2203.11171</u> , 2022b.
642 643 644	Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin, Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng He, Zilong Zheng, Yaodong Yang, et al. Jarvis-1: Open-world multi-task agents with memory-augmented multimodal language models. <u>arXiv preprint arXiv:2311.05997</u> , 2023c.
646 647	Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents. <u>arXiv</u> preprint arXiv:2302.01560, 2023d.

648 649 650	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in Neural Information Processing Systems, 35:24824–24837, 2022.
652 653 654	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents: A survey. <u>arXiv preprint arXiv:2309.07864</u> , 2023.
655 656 657	Jiannan Xiang, Tianhua Tao, Yi Gu, Tianmin Shu, Zirui Wang, Zichao Yang, and Zhiting Hu. Language models meet world models: Embodied experiences enhance language models. <u>Advances</u> <u>in neural information processing systems</u> , 36, 2024.
658 659 660 661	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. <u>arXiv preprint arXiv:2210.03629</u> , 2022.
662 663 664	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. <u>arXiv</u> preprint arXiv:2305.10601, 2023.
665 666 667	Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Keming Lu, Chuanqi Tan, Chang Zhou, and Jingren Zhou. Scaling relationship on learning mathematical reasoning with large language models. <u>arXiv preprint arXiv:2308.01825</u> , 2023.
669 670 671	Ceyao Zhang, Kaijie Yang, Siyi Hu, Zihao Wang, Guanghe Li, Yihang Sun, Cheng Zhang, Zhaowei Zhang, Anji Liu, Song-Chun Zhu, et al. Proagent: Building proactive cooperative ai with large language models. <u>arXiv preprint arXiv:2308.11339</u> , 2023.
671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 683 684 685 686 687 688 689 690 691 692 693 694 695 696 697 698 699	Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, et al. Ghost in the minecraft: Generally capable agents for open-world enviroments via large language models with text-based knowledge and memory. <u>arXiv preprint</u> <u>arXiv:2305.17144</u> , 2023.
700	

702 A METHOD DESCRIPTIONS

ExpertSystem: ExpertSystem simulates a human player's decision-making procedure. The approach includes a state evaluation function that calculates the battle advantage between the player's Pokémon and the opponent's, taking into account factors such as type-effectiveness, current HP and stats of Pokémon of both sides. If the current situation is favorable, the approach calculates the damage for attacks by considering move power, type effectiveness, and the stats of the Pokémon, and finally selects the move with the highest damage value; If the current situation is unfavorable and no

more advantageous situation using state evaluation function.
Overall, ExpertSystem is programmed with a decision-making procedure with numeric damage calculation, enabling it to generate effective actions for both move selection and switching, serving as a competitive opponent in our experiments.

strategic moves are available, it evaluates whether switching to another Pokémon could lead to a

RFT and DPO fine-tuning: We first supervisedly fine-tune LLaMA2-7B on 10,240 sampled frames (steps) of the expert system to ensure it 100% outputs admissible actions, as the initial LLM. We run the self-evolution iteration for 5 times, and in each iteration, we sample 100,000 frames (steps). The batch size for training is set to 64, and each sample is trained once. The learning rate is set to 1e-5 for RFT and 1e-6 for DPO based on empirical hyper-parameter tuning. RMSProp is adopted as the optimizer. For DPO, β is set to 0.1.

721 Reasoning approaches:

710

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738 739

740 741

742 743

744 745

746

- In-Out Prompt (IOPrompt): At each time step, LLM takes the observation as the input and directly output an action. We do not use few-shot examples as they do not provide obvious improvement (Table 3).
- Chain-of-Thoughts (CoT) (Wei et al., 2022): At each step, LLM generates thoughts to analyze the battle-field, then outputs an action conditioned on the thoughts. The reasoning includes the comparison of stats, type-effectiveness and evaluation of moves.
 - Self Consistency (SC-CoT) (Wang et al., 2022b): At each step, LLM independently generates *b* reasoning branches and do majority voting for the final output action.
 - Tree-of-Thoughts (ToT) (Yao et al., 2023): At each step, LLM analyzes the battlefield and proposes *b* top actions. LLM then criticizes and evaluates these actions with reasoning and selects the best action to output.
 - Reflexion (Shinn et al., 2023): LLM reflects on the outcome of the action a_{t-1} taken in the previous step, and use the reflection to generate action a_t in the next time step.
 - Last-Thoughts: At each step, LLM generates CoT reasoning by taking into consideration the thoughts from the last step.

B PROMPT DESCRIPTIONS

B.1 IOPROMPT

Input: observation o_t (game knowledge is presented in bold)

[System] You are playing a Pokemon battle and the goal is to win. Select the best action and output.

Historical turns:

748 Turn 12: Bouffalant used Megahorn, which was super effective to opposing Reuniclus. It 749 damaged opposing Reuniclus's HP by 59% (41% left). Opposing Reuniclus used Focus 750 Blast, which was super effective to Bouffalant. It damaged Bouffalant's HP by 53% (0% left). 751 Bouffalant fainted. Bouffalant outspeeded opposing Reuniclus. You sent out Doublade. 752 Current turn: 753 Opponent has 2 pokemons left. Opponent current pokemon:reuniclus, 754 Type:PSYCHIC,HP:41%,Atk:157,Def:174,Spa:258,Spd:191,Spe:99 755

757	reuniclus as defender, BUG,GHOST deal 2x damage; PSYCHIC,FIGHTING only deal
758	0.5x damage to reuniclus
759	reuniclus as attacker, PSYCHIC deal 2x damage to FIGHTING pokemon; PSYCHIC
760	deal 0.5X damage to PSY CHIC, STEEL pokemon Your current pokemon:doublade
761	Type: STEEL GHOST HP:10% Atk:228 Def:293 Spa:121 Spd:128 Spe:105
762	Your doublade has 4 moves can take:
763	swordsdance:Type:NORMAL,Cate:Status,Power:0,Acc:100%,
764	Effect:Raises the user's Attack by two stages.
765	ironhead:Type:STEEL,Cate:Physical,Power:105,Acc:100%,
766	Effect:Has a 30% chance to make the target flinch.
767	closecombat:Type:FIGHTING,Cate:Physical,Power:157,Acc:100%,
768	Effect: Lowers the user's Defense and Special Defense by one stage after inflicting damage.
769	shadowsneak: Type: GHOS 1, Cate: Physical, Power: 52, Acc: 100%,
770	Vou have 2 pokemons can switch:
771	azelf:Tyne:PSYCHIC HP:0% Atk:249 Def:160 Sna:249 Snd:160 Sne:233
772	Moves:[uturn.BUG].[nsvchic.PSYCHIC].[fireblast.FIRE]
773	tsareena:Type:GRASS.HP:100%.Atk:250.Def:213.Spa:132.Spd:213.Spe:169.
774	Moves:[rapidspin,NORMAL],[tripleaxel,ICE],[powerwhip,GRASS]
775	
776	Output: action a_t
777	
778	Action a.: Move: shadowsneak
779	Acton a_t . Move. shadowsheak
785 786	Input : observation o_t (game knowledge is presented in bold)
787 788 789 790 791	[System] You are playing a Pokemon battle and the goal is to win. Below is an example to teach you how to reason and select an action. ===Example Start===
791	===Example End===
793	Below is the real case:
794	Historical turns: Turn 20: Opposing Salaggla used Elemethrower which demaged Ducknein's UD by 120/ (00/
795	left) Dusknoir fainted. You sent out Lickilicky
796	Current turn:
797	Opponent has 1 pokemons left. Opponent current pokemon:salazzle,
798	Type:POISON&FIRE,HP:35%,Atk:152,Def:146,Spa:229,Spd:146,Spe:239
799	salazele as a defender, FIGHTING attack only deal 0.5x damage to it
	Your current pokemon:lickilicky.
800	,
800 801	Type:NORMAL,HP:44%,Atk:195,Def:213,Spa:187,Spd:213,Spe:135
800 801 802	Type:NORMAL,HP:44%,Atk:195,Def:213,Spa:187,Spd:213,Spe:135 Your lickilicky has 4 moves can take:
800 801 802 803	Type:NORMAL,HP:44%,Atk:195,Def:213,Spa:187,Spd:213,Spe:135 Your lickilicky has 4 moves can take: wish:Type:NORMAL,Cate:Status,Power:0,Acc:100%
800 801 802 803 804	Type:NORMAL,HP:44%,Atk:195,Def:213,Spa:187,Spd:213,Spe:135 Your lickilicky has 4 moves can take: wish:Type:NORMAL,Cate:Status,Power:0,Acc:100% Effect:User will recover half its max HP at the end of the next turn.
800 801 802 803 804 805	Type:NORMAL,HP:44%,Atk:195,Def:213,Spa:187,Spd:213,Spe:135 Your lickilicky has 4 moves can take: wish:Type:NORMAL,Cate:Status,Power:0,Acc:100% Effect:User will recover half its max HP at the end of the next turn. bodyslam:Type:NORMAL,Cate:Physical,Power:114,Acc:100%
800 801 802 803 804 805 806	Type:NORMAL,HP:44%,Atk:195,Def:213,Spa:187,Spd:213,Spe:135 Your lickilicky has 4 moves can take: wish:Type:NORMAL,Cate:Status,Power:0,Acc:100% Effect:User will recover half its max HP at the end of the next turn. bodyslam:Type:NORMAL,Cate:Physical,Power:114,Acc:100% Effect:Has a 30% chance to paralyze the target. knockoff:Ture:DAPK Cate:Physical Power:97 A act 100%
800 801 802 803 804 805 806 807	Type:NORMAL,HP:44%,Atk:195,Def:213,Spa:187,Spd:213,Spe:135 Your lickilicky has 4 moves can take: wish:Type:NORMAL,Cate:Status,Power:0,Acc:100% Effect:User will recover half its max HP at the end of the next turn. bodyslam:Type:NORMAL,Cate:Physical,Power:114,Acc:100% Effect:Has a 30% chance to paralyze the target. knockoff:Type:DARK,Cate:Physical,Power:87,Acc:100% Effect:Target drops its held item
800 801 802 803 804 805 806 807 808	Type:NORMAL,HP:44%,Atk:195,Def:213,Spa:187,Spd:213,Spe:135 Your lickilicky has 4 moves can take: wish:Type:NORMAL,Cate:Status,Power:0,Acc:100% Effect:User will recover half its max HP at the end of the next turn. bodyslam:Type:NORMAL,Cate:Physical,Power:114,Acc:100% Effect:Has a 30% chance to paralyze the target. knockoff:Type:DARK,Cate:Physical,Power:87,Acc:100% Effect:Target drops its held item. swordsdance:Type:NORMAL,Cate:Status Power:0 Acc:100%
800 801 802 803 804 805 806 807 808 809	Type:NORMAL,HP:44%,Atk:195,Def:213,Spa:187,Spd:213,Spe:135 Your lickilicky has 4 moves can take: wish:Type:NORMAL,Cate:Status,Power:0,Acc:100% Effect:User will recover half its max HP at the end of the next turn. bodyslam:Type:NORMAL,Cate:Physical,Power:114,Acc:100% Effect:Has a 30% chance to paralyze the target. knockoff:Type:DARK,Cate:Physical,Power:87,Acc:100% Effect:Target drops its held item. swordsdance:Type:NORMAL,Cate:Status,Power:0,Acc:100% Effect:Raises the user's Attack by two stages.

810	
811	You have 2 pokemons can switch: togekiss:Type:EAIRV&ELVING HD:28% Atk:85 Def:108 Spa:238 Spd:230 Spe:174
812	Moves [airs]ash El VING] [aurasnhere EIGHTING]
813	wobbuffet:Type:PSYCHIC.HP:48%.Atk:69.Def:167.Spa:119.Spd:167.
814	
815	Dutput:
817	
818	Thoughts: The opponent's Salazzle is faster than my Lickilicky, and it has a high Special
819	Attack stat. Lickilicky's current HP is 44%, which means it might not survive a strong special
820	attack from Salazzle. Switching to Togekiss is risky due to its low HP (28%), and Wobbuffet
821	might not be able to counter effectively due to its lower defenses. Lickilicky's best move is to
822	use Knock Off, which can deal damage and potentially remove Salazzle's held item, reducing
823	damage before fainting
824	Action a_t : Move: knockoff
825	
826	
827 I	3.3 TREE-OF-THOUGHTS (TOT)
828 T	For each step. ToT has two phases: proposal and self-evaluation
829	or each step, for has two phases. proposal and sen-evaluation.
830	Phase 1: Proposal
⁸³¹ I	input: observation o_t (game knowledge is presented in bold)
832	
033	[System] You are playing a Pokemon battle and the goal is to win. Below is an example to
835	teach you how to reason and propose 3 best actions.
836	===Example Start===
837	
838	Here is the real case:
839	Historical turns: Turn 3: opposing Thievul used Dark Pulse. It damaged Dragalge's HP by 78%
840	(22% left). Dragalge used Focus Blast. It missed. opposing Thievul outspeeded Dragalge.
841	Current turn:
842	Opponent has 6 pokemons left.
843	Opponent current pokemon: thievul,
844	Type:DAKK, HP:52%, Alk:152, Del:152, Spa:009(4 stage), Spa:212, Spe:209 thiowyl as defender, FICHTING deal 2x damage: CHOST DARK only deal 0.5x damage:
845	PSVCHIC have no effect to thievul
846	thievul as attacker, DARK deal 2x damage to PSYCHIC, GHOST pokemon; DARK deal
847	0.5x damage to DARK, FIGHTING pokemon
848	Your current pokemon:dragalge,
849	Type:POISON&DRAGON,HP:28%,Atk:135,Def:206,Spa:219,Spd:264,Spe:126
850	Your dragalge has 4 moves can take:
851	Effect. Scatters noisoned snikes noisoning apposing Pakémon that switch in
852	dracometeor: Type: DRAGON. Cate: Special Power: 134. Acc: 90%.
000	Effect: Lowers the user's Special Attack by two stages after inflicting damage.
004 955	sludgewave:Type:POISON,Cate:Special,Power:98,Acc:100%,
856	Effect:Has a 10% chance to poison the target.
857	focusblast:Type:FIGHTING,Cate:Special,Power:124,Acc:70%,
858	Effect: Has a 10% chance to lower the target's Special Defense by one stage.
859	YOU HAVE 3 POKEMONS CAN SWITCH: blacephalon:Type:FIREGHOST HP:100% Atta: 205 Def:120 Spec224 Spd:170 Spec215
860	Moves: [fireblast FIRE] [shadowhall GHOST] [psychock PSYCHIC]
861	jellicent:Type:WATERGHOST,HP:100%.Atk:107.Def:168.Sna:193.Snd:227.Sne:151
862	Moves:[shadowball,GHOST],[scald,WATER]
863	snorlax:Type:NORMAL,HP:100%,Atk:228,Def:154,Spa:154,Spd:228,Spe:96,

I

Moves:[doubleedge,NORMAL],[darkestlariat,DARK],[facade,NORMAL] slowbro:Type:WATERPSYCHIC,HP:100%,Atk:132,Def:236,Spa:219,Spd:185,Spe:100, Moves:[scald,WATER] regieleki:Type:ELECTRIC,HP:100%,Atk:199,Def:122,Spa:199,Spd:122,Spe:353, Moves:[explosion,NORMAL],[thunderbolt,ELECTRIC],[voltswitch,ELECTRIC] Output:

Thoughts: Thievul's previous Dark Pulse did 78% damage, leaving Dragalge at 22% HP. Since Thievul outspeeds Dragalge, it will KO Dragalge before it can attack. Switching out Dragalge keeps it as fodder or for a potential safe switch-in later. If choose to switch, the switch in pokemon will first bear an attack. With a high SpD, Snorlax can potentially survive a boosted Dark Pulse from Thievul. If Thievul misses, we can use focusblast to knock out Thievul. **Proposal 1**: Switch: snorlax **Proposal 2**: Move: focusblast

Proposal 3: switch: dracometeor

Phase 2: Self-evaluation

Input: The input is the observation o_t and the proposal in phase 1, along with an example to teach LLM for self-evaluation. We omit the prompt here due to similarities in content.

Ouput:

864

865

866

867

868

870

871 872

873

874

875

876

877

878

879

880

882

883

884 885

886 887

889

890

891

892

893 894 895

896 897

202

Thoughts: Since Thievul outspeeds Dragalge and has high Spa, it will KO Dragalge before it can attack. If choose to use attack, it is very likely be knocked out in the next step. However, if choose to switch in a non force switch, the switch in pokemon will first bear an attack and be knocked out since Thievul's Special Attack has been boosted very high. Therefore, it is better to let Thievul stay in the battle to be knocked out and switch in a faster pokemon in a force switch.

Action a_t : Move: focusblast

B.4 LAST-THOUGHTS

Input: observation o_t (last thoughts, game knowledge are presented in **bold**)

899	[System] You are playing a Pokemon battle and the goal is to win. Below is an example to
900	teach you how to reason and select an action
901	===Example Start===
902	
903	===Example End===
904	Here is the real case:
905	Historical turns:
906	Turn 2: Froslass used Poltergeist, which damaged opposing Cobalion's HP by 29% (30%
907	left). Opposing Cobalion used Iron Head, which was super effective to Froslass and damaged
908	Froslass's HP by 100% (0% left). Froslass fainted. Froslass outspeeded opposing Cobalion.
909	You sent out Darmanitan.
910	Current turn:
911	Opponent has 6 pokemons left.
912	Opponent current pokemon: cobalion
913	Type:STEELFIGHTING, HP:30%, Alk:570(2 stage), De1:240, Spa:185, Spa:157, Spe:214
914	contained as detended, FIGHTING, GROUND, FIRE deal 2x damage; ICE, GRASS, DARK only deal 0.5x damage: ROCK BUC only deal 0.25x damage to cobalian
915	cobalion as attacker STEFI deal 2x damage to ICE ROCK pokemon: STEFI deal 0.5x damage
016	to FIRE pokemon: FIGHTING deal 2x damage to ICE ROCK DARK pokemon: FIGHTING
017	deal 0.5x damage to PSYCHIC BUG pokemon
911	

918	
919	Your current pokemon:darmanitangalar,
920	Type:ICE,HP:100%,Atk:263,Def:131,Spa:92,Spd:131,Spe:193
921	Your darmanitangalar has 4 moves can take:
922	flareblitz:Type:FIRE,Cate:Physical,Power:128,Acc:100%,
923	Effect: User takes 1/3 the damage inflicted in recoil. Has a 10% chance to burn the target.
924	uturn: Type:BUG,Cate:Physical,Power: 75,Acc:100%,
925	Effect: User must switch out after attacking.
026	earthquake: Type: GROUND, Cate: Physical, Power: 107, Acc: 100%,
920	Effect: Inflicts regular damage and can hit Dig users.
927	Effects Here a 2007 shores to make the torget flingh
928	Effect: Has a 50% chance to make the target minch.
929	crustle: Type: BLIGBOCK HD: 100% Att. 222 Def: 255 Spa: 156 Spd: 172 Spe: 122
930	Moves:[stoneedge ROCK] [earthquake GROUND] [vscissor BUG]
931	grapploct Type FIGHTING HP 100% Atk 252 Def 204 Spa 170 Spd 187 Spe 121
932	Moves ficepunch ICE] [drainpunch FIGHTING] [suckerpunch DARK]
933	audino: Type: NORMAL HP: 100%. Atk: 161. Def: 208. Spa: 161. Spd: 208. Spe: 143.
934	Moves:[knockoff.DARK]
935	mesprit:Type:PSYCHIC,HP:100%,Atk:181,Def:225,Spa:225,Spd:225,Spe:183,
936	Moves: [psychic, PSYCHIC], [icebeam, ICE], [energyball, GRASS]
937	Your thought from last step: Cobalion has boosted its Attack by 2 stages, making it a signifi-
938	cant threat. It is weak to FIGHTING, GROUND, and FIRE moves. Crustle has Earthquake, a
939	GROUND move that is super effective against Cobalion, but Crustle is slower and vulnerable to
940	Cobalion's FIGHTING moves. Grapploct has Drain Punch, a FIGHTING move that is super ef-
941	fective, but it is also slower and vulnerable to Cobalion's FIGHTING moves. Darmanitan-Galar
942	has Flare Blitz, a FIRE move that is super effective, and it is faster than Cobalion, allowing it to
943	attack first. Audino and Mesprit do not have super-effective moves and are not ideal choices.
944	Darmanitan-Galar is the best choice due to its speed advantage and super-effective Flare Blitz.
945	
946	Output:
0.47	-

Thoughts: Darmanitan-Galar is currently on the field and has a speed advantage over Cobalion, allowing it to attack first. Cobalion is at 30% HP and is weak to FIRE moves. Darmanitan-Galar's Flare Blitz is a FIRE move that will deal 2x damage to Cobalion. Although Flare Blitz has recoil damage, the priority is to eliminate the immediate threat posed by Cobalion's boosted Attack. Therefore, using Flare Blitz is the best option to potentially knock out Cobalion before it can attack.

Action a_t : Move: flareblitz

We can observe from the feedback described in o_t that our Froslass fainted in the last turn, triggering a force switch. The LLM chose to switch in Darmanitan with the thoughts from the last step and opted to attack with Flare Blitz in this step. If the last thought is not provided, given that the opponent has already doubled its Attack stats, the LLM is very likely to behave inconsistently, as described in Appendix C.

B.5 REFLEXION

948

949

950

951

952

953

954 955

956

957

958

959 960 961

962

963 964

965

97 97

Phase 1: Reflection on the last step t - 1

Input: Observation o_{t-1} , Action a_{t-1} , Feedback f_{t-1}

966	[System] Make a reflection given the state action and its outcome in a Pokémon hattle Below
967	is an example to teach you how to reflect on a previous battle step.
968	===Example Start===
969	
970	===Example End===
971	Here is the real case:

070	
972	Historical turns: Turn 12: opposing Sylveon used Psyshock. It was super effective to Dragalge.
973	It damaged Dragalge's HP by 56% (0% left). Dragalge fainted. You sent out Clefable.
975	Current turn:
976	Opponent has 4 pokemons left.
977	Opponent's known pokemon off the field:drednaw
978	Opponent current pokemon:sylveon, The FAIDVIID 100% Add 156 D 6 156 See 245(1,4444) See 1204(1,4444) See 147
979	Type:FAIRY, HP:100%, AIK:150, De1:150, Spa:545(1 stage), Spa:594(1 stage), Spe:14/
980	effect to sylveon
981	sviveon as attacker. FAIRY deal 2x damage to FIGHTING.DRAGON.DARK pokemon:
982	FAIRY deal 0.5x damage to FIRE pokemon
983	Your current pokemon:clefable,
984	Type:FAIRY,HP:100%,Atk:119,Def:167,Spa:203,Spd:195,Spe:146 Your clefable has 4 moves
985	can take:
986	moonblast:Type:FAIRY,Cate:Special,Power:73,Acc:100%,
987	Effect: Has a 30% chance to lower the target's Special Attack by one stage.
988	calmmind: Type:PSYCHIC, Cate:Status, Power:0, Acc: 100%,
989	fireblast Type FIRE Cate Special Power: 85 Acc: 85%
990	Effect:Has a 10% chance to burn the target.
991	softboiled:Type:NORMAL.Cate:Status.Power:0.Acc:100%.
992	Effect:Heals the user by half its max HP.
993	You have 4 pokemons can switch:
994	galvantula:Type:BUGELECTRIC,HP:19%,Atk:131,Def:146,Spa:206,Spd:146,Spe:224,
995	Moves:[thunder,ELECTRIC],[bugbuzz,BUG],[voltswitch,ELECTRIC]
996	turtonator:Type:FIREDRAGON,HP:39%,Atk:136,Def:275,Spa:201,Spd:191,Spe:109,
997	Moves: [fireblast, FIRE], [dracometeor, DRAGON]
998	Moves:[superpower FIGHTING] [knockoff DARK]
999	zarudedada:Type:DARKGRASS HP:63% Atk:232 Def:209 Spa:154 Spd:193 Spe:209
1000	Moves:[powerwhip.GRASS].[darkestlariat.DARK]
1001	Chosen Action (a_{t-1}) : Move: fireblast
1002	Outcome (f_{t-1}) : Clefable used Fire Blast, which damaged opposing Sylveon's HP by 16%
1003	(84% left). Opposing Sylveon outspeeded Clefable.
1004	
1005	Output: Reflection
1006	
1007	Reflection : The chosen action to use Fire Blast against Sylveon was not optimal. Sylveon, being
1008	a Fairy-type, does not have a weakness to Fire-type moves, and the damage dealt was minimal.
1009	Additionally, Sylveon outspeeded Clefable, allowing it to attack first and deal significant damage
1010	with Hyper Voice, reducing Clefable's HP by 65%. Given Sylveon's high Special Defense and
1011	the fact that it had a Special Attack boost, a better strategy would have been to use Moonblast,
1012	which could potentially lower Sylveon's Special Attack, or to use Calm Mind to boost Clerable's
1013	special Defense and Special Attack, preparing for a more effective counterattack. Switching
1014	been considered, although its low HP would have made it a risky move
1015	oven considered, annough its low iti would have made it a lisky move.
1016	

1017 Phase 2: Action generation

Input: Observation o_t , Reflection on the last step

1020
1021
1021
1021
1022
1022
1023
1024
1024
1024
1025
1025
1026
1027
1028
1029
1029
1029
1020
1020
1021
1021
1021
1022
1023
1024
1025
1025
1024
1025
1025
1025
1026
1027
1027
1028
1029
1029
1020
1020
1020
1021
1021
1021
1021
1022
1023
1024
1025
1024
1025
1025
1025
1026
1027
1027
1028
1028
1029
1029
1020
1020
1020
1020
1020
1021
1021
1021
1021
1022
1022
1022
1023
1024
1025
1025
1025
1024
1025
1025
1025
1025
1025
1025
1026
1027
1026
1027
1028
1028
1029
1029
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020
1020</l

1026	
1027	Opponent has 4 pokemons left. Opponent's known pokemon off the field:drednaw
1028	Opponent current pokemon:sylveon,
1029	Type:FAIRY,HP:90%,Atk:156,Def:156,Spa:345(1 stage),Spd:394(1 stage),Spe:147
1030	sylveon as defender, FIGHTING, DARK, BUG only deal 0.5x damage; DRAGON have no
1031	ellect to sylveon sylveon as attacker FAIDV deal 3x damage to FICHTINC DDACON DADK nekement
1032	FAIRY deal 0.5v damage to FIRF nokemon
1033	Your current nokemon: clefable Type: FAIRY HP:35% Atk: 119 Def:167 Spa: 203 Spd:195 Spe:146
1034	Your clefable has 4 moves can take:
1035	moonblast:Type:FAIRY.Cate:Special.Power:73.Acc:100%.
1036	Effect: Has a 30% chance to lower the target's Special Attack by one stage.
1037	calmmind:Type:PSYCHIC,Cate:Status,Power:0,Acc:100%,
1038	Effect:Raises the user's Special Attack and Special Defense by one stage.
1039	fireblast:Type:FIRE,Cate:Special,Power:85,Acc:85%,
1040	Effect:Has a 10% chance to burn the target.
1041	softboiled:Type:NORMAL,Cate:Status,Power:0,Acc:100%,
1042	Effect: Heals the user by half its max HP.
1043	You have 4 pokemons can switch:
1044	galvaniula: Type: DUGELECTRIC, HP: 19%, Alk: 151, Del: 140, Spa: 200, Spu: 140, Spe: 224, Moves: [thunder ELECTPIC] [hughuzz BUG] [voltswitch ELECTPIC]
1045	turtonator: Type: FIREDRAGON HP:39% Atk:136 Def:275 Sna:201 Snd:191 Sne:109
1046	Moves: [fireblast FIRE] [dracometeor DR AGON]
1047	malamar:Type:DARK&PSYCHIC.HP:100%.Atk:193.Def:187.Spa:155.Spd:166.Spe:163.
1048	Moves:[superpower,FIGHTING],[knockoff,DARK]
1049	zarudedada:Type:DARK&GRASS,HP:63%,Atk:232,Def:209,Spa:154,Spd:193,Spe:209,
1050	Moves:[powerwhip,GRASS],[darkestlariat,DARK]
1051	Refection on the last step : The chosen action to use Fire Blast against Sylveon was not optimal.
1052	Sylveon, being a Fairy-type, does not have a weakness to Fire-type moves, and the damage
1053	dealt was minimal. Additionally, Sylveon outspeeded Clefable, allowing it to attack first and
1054	deal significant damage with Hyper Voice, reducing Clefable's HP by 65%. Given Sylveon's
1055	high Special Delense and the fact that it had a Special Attack boost, a better strategy would have been to use Meenblest, which could potentially lower Sylveon's Special Attack or to
1056	use Calm Mind to boost Clefable's Special Defense and Special Attack, preparing for a more
1057	effective counterattack. Switching to a Pokémon with higher resistance to Fairy-type moves
1058	such as Turtonator, could have also been considered, although its low HP would have made it a
1059	risky move.
1060	

Output:

Action a_t : Switch: malamar

We can observe that the LLM is not satisfied with the action in the last step and choose to switch in current step, increasing the inconsistency.

C ACTION INCONSISTENCY



Figure 5: When facing a powerful Pokémon, the LLM switches different Pokémon in three consecutive steps to prevent its Pokémon from fainting. However, this gives the opponent three free turns to quadruple its attack stats and quickly defeat the agent's entire team.

1080 TYPE 1082 1084 1085 1086 1087 1088 1089 1090 1091 1093 1094 Figure 6: Type effectivness chart. "+" denotes super-effective (2x damage); "-" denotes ineffective 1095 (0.5x damage); "×" denotes no effect (0x damage). Unmarked is standard (1x) damage. 1096 1099 Species: Charizard Move: 1100 Dragon Dance Ability: Blaze Fire Blast [120] 1101 HP: 319 Atk: 293 Outrage [120] 1102 Def: 280 Spe: 328 [75] Air Slash 1103 GRASS (POISON) Species: Venusaur Move: 1104 Synthesis Ability: Overgrow 1105 Solar Beam [120] GRASS HP: 364 Atk: 289 1106 Sludge Bomb [90] (POISON Def: 328 Spe: 284 [75] GRASS Giga Drain 1107 1108 Figure 7: Two representative Pokémon: Charizard and Venusaur. Each Pokémon has type(s), ability, 1109 stats and four battle moves. 1110 1111 1112 1113 We observe that when facing a powerful Pokémon, LLMs are more tend to behave inconsistently. 1114 As illustrated in Figure 5, starting from turn 8, the agent chooses to continuously switch to different 1115 Pokémon in three consecutive turns, giving the opposing Pokémon three free turns to boost its 1116 attack stats to four times and take down the agent's entire team quickly. This phenomenon can be 1117 exacerbated by CoT but addressed by Last-Thoughts. 1118 1119 1120 POKÉMON KNOWLEDGE D 1121 1122 Species: There are more than 1,000 Pokémon species (bul, 2024b), each with its unique ability, 1123 type(s), statistics (stats) and battle moves. Figure 7 shows two representative Pokémon: Charizard 1124 and Venusaur. 1125

Type: Each Pokémon species has up to two elemental types, which determine its advantages and weaknesses. Figure 6 is the type-effectiveness chart that presents relationship between 18 types of attack moves and attacked Pokémon. For example, 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, 2024a).