Pokémon Red via Reinforcement Learning

Marco Pleines*, Daniel Addis[‡], David Rubinstein[†],
Frank Zimmer¶, Mike Preuss∥, Peter Whidden[†]
*TU Dortmund University, Dortmund, Germany

†Independent Researcher, Brooklyn, NY, USA, ‡Independent Researcher, Minneapolis, MN, USA

¶Rhine-Waal University of Applied Sciences, Kamp-Linfort, Germany

∥LIACS Universiteit Leiden, Leiden, Netherlands

Abstract—We present a Deep Reinforcement Learning (DRL) agent that successfully completes the first several hours of Pokémon Red, a classic Game Boy JRPG that exposes significant challenges as a testbed for agents, including multitasking, long horizons of tens of thousands of steps, hard exploration, and a vast array of potential policies. Our agent completes an initial segment of the game, up to Cerulean City, a location requiring progression through two cities, battle-filled passages, a maze-like cave, and defeating the first gym leader. Our experiments include various ablations that reveal vulnerabilities in reward shaping. We argue that long-form games like Pokémon hold strong potential for future research, presenting long-term coherent reasoning challenges absent from simpler arcade games. Our environment wrapper, training algorithm, human replay data, and pretrained agent are available at (REDACTED FOR REVIEW).

Index Terms—pokemon, deep reinforcement learning, reward shaping, proximal policy optimization

I. Introduction

Video games have significantly advanced artificial intelligence, acting as both a catalyst for research and a benchmark for solving complex problems. DRL and its adjacent fields have achieved remarkable results by enabling the training of agents to play games with extremely long horizons such as DotA 2 [1], NetHack [2], and Minecraft [3]. While these modern games were used to demonstrate strong advancements, older games remain an untapped resource, offering similar levels of complexity, unique challenges, and opportunities for research. This work examines an older game, the 1996 Game Boy JRPG: Pokémon Red.

Pokémon games have consistently gained widespread attention, both within and beyond the gaming community. A 2023 YouTube video showcasing an early DRL agent playing Pokémon Red went viral, garnering over 7.5 million views [4]. The video's viral success highlights the game's broad appeal and serves as the inspiration that we build upon. Pokémon and its many successors challenge the player with a variety of tasks with the ultimate goal of becoming a Pokémon Master by defeating the game's Champion and catching all 151 Pokémon species. In this paper, we highlight Pokémon Red as a game that exposes players to multiple complex tasks:

- Multi-Task Challenge: Pokémon includes strategic battling, 2D navigation, UI navigation, resource management, and puzzle-solving (Figure 1).
- Exploration Challenge: Beating Pokémon requires careful planning and navigation. In our experiments, the trained

- DRL agents made tens of thousands of decisions per episode to complete only a small fraction of the game.
- Large Policy Space: With 151 Pokémon species, diverse combinations of moves, and numerous usable items, Pokémon Red features vast gameplay possibilities and freedom to develop successful policies.

To explore Pokémon Red's potential, we contribute a minimal DRL environment and a simple, adaptable baseline agent trained with Proximal Policy Optimization (PPO) [5]. Our approach lays the groundwork for more advanced exploration and hierarchical methods, while offering detailed insights into the performance and behavior of the trained baseline agent and various reward ablations. This paper is structured as follows: we establish the Pokémon baseline by introducing the gameplay and properties of the Pokémon Red environment; then, the training setup is detailed. Section IV presents results, discussing the key findings and potential future directions. Before concluding, we refer to related work.

II. POKÉMON RED ENVIRONMENT

This section describes the core mechanics of Pokémon Red and our formulation as a Markov Decision Process (MDP) [6], including the design of the observation space, action space, reward function, and terminal conditions. The simulation speed of the environment was measured using random actions on an AMD Ryzen 7 2700X CPU. Pokémon runs at ≈9403 steps per second (SPS), which is faster than Atari Breakout (7117 SPS) and slower than Proceen (18530 SPS) [7].

A. Selected Game Objectives

In this work, we will focus on Pokémon Red's storyline objectives, which entail completing several quests and milestones. For simplicity, we have reduced the objective from completing the entire game to completing only the main tasks up to the end of Cerulean City, the location of the second gym. This constitutes approximately 20% of the full game and is achieved by accomplishing the following milestones:

- Start in Pallet Town and pick a starter Pokémon from Professor Oak.
- Travel to Viridian City and obtain Oak's Parcel.
- Return the parcel to Oak's lab in Pallet Town.
- Traverse Viridian Forest to reach Pewter City.
- Defeat Pewter City's gym leader, Brock (Badge 1).
- Traverse Mt. Moon to arrive at Cerulean City.

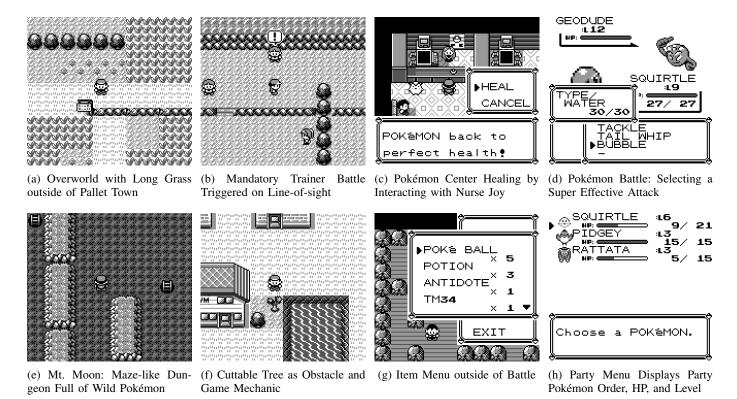


Fig. 1: All subfigures depict game screens possibly perceivable by the agent, presenting various challenges in Pokémon Red that involve exploration, navigation, and strategic decision-making. The agent must traverse a 2D overworld (a) with a party of Pokémon (h), winning mandatory trainer battles (b) and navigating complex maze-like areas (e). Progression also depends on interacting with the game's UI, including healing Pokémon at a Pokémon Center (c), using strategy to win battles (d), using items (g). Overcoming obstacles, like cuttable trees (f), requires the use of mandatory game-mechanic moves, which must be obtained via exploration, taught to eligible Pokémon, and used via the Start menu interface when facing the obstacle.

- Defeat Cerulean City's gym leader, Misty (Badge 2).
- Transform Bill, a storyline non-playable character (NPC), back into a human to unlock the path to Vermilion City.

Because the first two objectives require backtracking, we have started the environment after the third objective, as the remaining tasks involve minimal or no backtracking. Thus, the agent does not choose its starter. Instead, we let it start with the water-type Pokémon Squirtle, which has an advantage against gym leader Brock's rock Pokémon. After traversing Mt. Moon, the subsequent storyline objectives can be completed in a variety of orders, and the game transitions into a large, openworld setting with numerous non-linear tasks and puzzles.

B. Game Mechanics

Pokémon Red features turn-based combat, strategic teambuilding, and exploration of an expansive grid world. The player navigates interconnected maps, overcoming obstacles such as trainer battles, environmental barriers, and puzzles. Early on, battles are the primary challenge.

Battles initiate with the first Pokémon in each player's party (a wild Pokémon is treated as a party of one). Battles are turn-based, with each side selecting moves that are executed in order based on the Pokémon's Speed stat, so that faster Pokémon act before slower ones. On each turn, a player may

choose a move (Fig. 1d), switch to another Pokémon in the party (Fig. 1h), use an item, or run from battle. Escaping from battle is only possible during wild encounters. Wild Pokémon appear randomly in grass, on water, or inside dungeons, and players can capture them using Poké Balls. Each Pokémon can have up to four moves, each with a finite number of Power Points (PP) representing the number of times a move can be used in battle. Using a move consumes one PP and may heal, deal damage, inflict status effects, or temporarily alter stats. The effects of a move can vary based on the attacking and target Pokémon's types. This type match-up system is analogous to rock-paper-scissors, where each type holds a natural advantage over another. For example, water-type moves deal double damage against fire-type Pokémon. Attacks can also randomly become Critical Hits that deal double damage. When a Pokémon's Hit Points (HP) reach zero, it faints and can no longer be used in battle. A Pokémon's HP and PP can be restored using items (Fig. 1g) or by visiting a Pokémon Center (Fig. 1c). Defeating a wild or trainer Pokémon grants experience points (XP) to all conscious Pokémon that participated in battle. To gain XP, a Pokémon must remain conscious until the enemy Pokémon faints. Gaining XP leads to levelups that increase a Pokémon's stats and offer opportunities to

learn new moves or evolve. Evolution is triggered at certain levels and can be canceled with the B button; it raises a Pokémon's stats but may delay learning new moves, offering a strategic tradeoff. Many Pokémon have evolved forms, and some moves are learned at lower levels in their base forms than after evolving. If the player's entire party faints, it is considered a "blackout," and the player will respawn at the last visited Pokémon Center or at their home in Pallet Town if no Pokémon Centers have been visited.

C. Observation Space

The agent receives two observation modalities: a vision input and a game state vector.

The vision input uses a grayscale representation of the Game Boy's screen downsampled by two for a resolution of 72×80 pixels. We stack the current frame and two previously observed frames to provide minimal temporal awareness. The agent additionally observes a separate 48×48 binary crop of the screen, centered on the player, that indicates visited coordinates.

The game state vector includes the current HP and level of each Pokémon in the party, as well as flags indicating the completion status of various in-game events. Notably, this observation space is intentionally limited for simplicity and to mimick the information a first time human player would have. The observation does not provide information on the species, stats, nor movesets of the Pokémon.

D. Action Space

The action space is discrete and consists of seven Game Boy buttons: A, B, Start, Up, Down, Left, and Right (Select is omitted). An agent must select one of these actions for each environment step. Navigation of both the overworld and in-game menus is performed using the arrow keys. The A button confirms selections or initiates an interaction, while the B button cancels actions or exits menus. The Start button opens the main menu.

For simplicity, actions are represented as full presses instead of having a press/unpress action for each button. To reliably move one tile at a time in the grid world, each button is held for 8 frames, then released for 16 frames, resulting in one decision every 24 frames. As a result, the simulation runs at a slower speed of approximately $\frac{9403}{24} \approx 392$, SPS.

E. Determinism and Terminal Conditions

Wild encounters, Pokémon stats, critical hits, and other mechanics appear stochastic. However, the random number generator (RNG) on the Game Boy is influenced by the player's input. As a result, if the player behaves deterministically, the environment will also behave deterministically. This property is commonly exploited by speedrunners and could potentially be leveraged by trained agents as well [4]. The agent could execute a specific sequence of inputs to encounter a strong Pokémon and catch it on the first try. To counteract this determinism, a sequence of random button presses could be executed upon resetting the environment.

Regarding terminal conditions, time is the only constraint we define. The agent begins with a budget of 10,240 steps, which increases by 2,048 steps for each completed event. Events include important trainer battles, distinct NPC interactions, and storyline progression. This dynamic step budget introduces variety in training samples, preventing environments from resetting at fixed intervals and mitigating the risk of catastrophic forgetting (Section IV-B2).

F. Reward Function

Pokémon has an inherently long time horizon, requiring hours of real-world play time before major objectives are achieved. We use dense auxiliary rewards to reinforce intermediate progress towards and behaviors that support reaching these goals.

Event Reward: The agent receives a reward of R_{event} = +2 for every completed event (trainer battle, step in quest progression).

Navigation Reward: Without a navigation reward, the agent converges to a policy that does not progress much further than the player's origin in Pallet Town. To aid in exploring the overworld, the agent receives R_{nav} = +0.005 for each new overworld coordinate visited within an episode.

Healing Reward: When HP is gained, the agent receives a reward proportional to the fractional sum of party HP:

$$R_{\text{heal}} = 2.5 \sum_{i=1}^{6} \frac{\text{HP}_{i}^{\text{after}} - \text{HP}_{i}^{\text{before}}}{\text{HP}_{i}^{\text{max}}} \tag{1}$$

This reward is triggered when a Pokémon levels up, when a Pokémon Center is used, or when a potion is used on a Pokémon with less than full HP; each of these behaviors supports progression.

Level Reward: The agent receives a reward based on the levels the Pokémon in its party. This encourages early level progression and the capture of new Pokémon while discouraging wild-encounter grinding. To balance this behavior, we define:

$$R_{\text{lvl}} = 0.5 \min\left(\sum_{i=1}^{6} \text{lvl}_i, \frac{\sum_{i=1}^{6} \text{lvl}_i - 22}{4} + 22\right)$$
 (2)

The threshold of 22 signifies a reasonable level to battle the leader of Gym 2, Misty, in Cerulean City, with the marginal gain per additional level downscaled by a factor of 4.

Finally, all rewards are summed up to yield $R = R_{\rm event} + R_{\rm nav} + R_{\rm heal} + R_{\rm lvl}$. Rewards are further discussed in Section IV-B5.

III. TRAINING METHOD

All agents are trained using Proximal Policy Optimization (PPO) [5]. We leverage the clipped surrogate loss to optimize

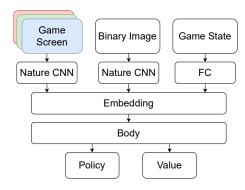


Fig. 2: Actor-Critic Network (2M parameters, GRU: 4M).

the policy at time t:

$$\begin{split} L_t^C(\theta) &= -\mathbb{E}_t \Bigg[\min \Big(q_t(\theta) A_\pi^{\text{GAE}}(s_t, a_t), \\ & \text{clip}(q_t(\theta), 1 - \epsilon, 1 + \epsilon) A_\pi^{\text{GAE}}(s_t, a_t) \Big) \Bigg] \\ \text{where} \quad q_t(\theta) &= \frac{\pi(a_t | s_t, \theta)}{\pi(a_t | s_t, \theta_{\text{old}})} \end{split} \tag{3}$$

 s_t denotes the observed state, a_t is the chosen action, θ defines the parameters of the policy, ϵ is the clip range, and $A_\pi^{\rm GAE}$ is the generalized advantage estimate. The value function shares parameters with the policy and undergoes clipping:

$$\begin{split} L_t^{VClip}(\theta) &= \Big(\mathrm{clip} \big(V(s_t, \theta), \\ &\quad V(s_t, \theta_{\mathrm{old}}) - \epsilon, \\ &\quad V(s_t, \theta_{\mathrm{old}}) + \epsilon \big) - \hat{V}_t \Big)^2 \end{split} \tag{4}$$

$$L_t^V(\theta) = \max \left[\left(V(s_t, \theta) - \hat{V}_t \right)^2, \ L_t^{VClip}(\theta) \right]$$
 (5)

The final combined loss is depicted by $L_t^{C+V}(\theta)$:

$$L_t^{C+V}(\theta) = \mathbb{E}_t[L_t^C(\theta) + vL_t^V(\theta)] \tag{6}$$

where v is the value function's coefficient.

We intentionally leave out the commonly used entropy bonus, as tuning its weighting across 10^{-3} , 10^{-4} , 10^{-5} , and 10^{-6} yielded no significant improvements.

A. Network Architecture

The general network architecture of the agent is illustrated in Figure 2. Each vision modality within the observation space is encoded using a classic Nature Convolutional Neural Network (CNN) [8]. The game state vector is processed separately through a single fully connected layer. The outputs of all encoders are then flattened, concatenated, and projected to match the dimensionality of the network's body, which consists of a fully connected layer as standard, or, to facilitate recurrence, a Gated Recurrent Unit (GRU). Pleines et al.

TABLE I: The hyperparameters used in our experiments.

Hyperparameter	Value	Hyperparameter	Value
Discount Factor γ	0.997	Clip Range ϵ	0.2
λ (GAE)	0.95	Value Function Coef. v	0.5
Number of Workers	32	Learning Rate	0.0003
Worker Steps	2048	Max Gradient Norm	0.5
Epochs	3	Activations	ReLU
Mini Batches	8	Optimizer	AdamW

[7] describe how the GRU enables memory capabilities by maintaining a hidden state over time. The hidden features computed in the body are used to construct both a policy head and a value head, which respectively output the discrete policy (actions) and state value estimate (via the value function). The ReLU activation function is applied throughout the network.

B. Hyperparameters

Table I presents the hyperparameters used across all experiments. Since Pokémon Red features episodes lasting tens of thousands of steps, each environment instance (i.e., worker) samples a horizon of 2048 steps, which we found to be more effective than horizons of 512 or 1024 steps. This longer horizon increases the likelihood of capturing more reward signals compared to the more commonly used shorter horizons, such as 512 steps or fewer. To leverage all available cores of our training hardware, we use 32 workers, resulting in a total batch size of 65,536. γ was tuned over $\{0.99, 0.997, 0.999, 0.9995\}$.

IV. EXPERIMENTAL RESULTS AND DISCUSSION

The previous sections introduced the MDP behind Pokémon Red and the training algorithm. These establish our minimal baseline, which we analyze through ablation experiments. We vary the agent's starting state by selecting Squirtle, Charmander, or Bulbasaur as the starter Pokémon. Additional experiments ablate individual reward signals and compare the results to 30 playthroughs by skilled human players.

Each experiment is repeated five times with different seeds. During training, we evaluate 25 time points using 30 episodes per repetition. We report the mean performance across runs and measure variability with the standard deviation.

A. Baseline Performance

Figure 3 shows sample efficiency curves, depicting the mean proportion of milestones completed. The experiments vary by starter Pokémon, an agent using memory via GRU, and a *Fast* setting with increased text speed and disabled battle animations.

The first milestone is reaching Mt. Moon after defeating Brock. *Fast* performs best, peaking at 98% completion rate, followed by the baseline at 97%. With Bulbasaur, the agent reaches 94% but fails entirely at the second milestone: arriving in Cerulean City. This outcome is due to exploiting the heal reward as discussed in Section IV-B3. Other agents achieve between 53% and 93%.

Beating Misty and completing Bill's quest can be done in any order, but only the latter unlocks Vermilion City.

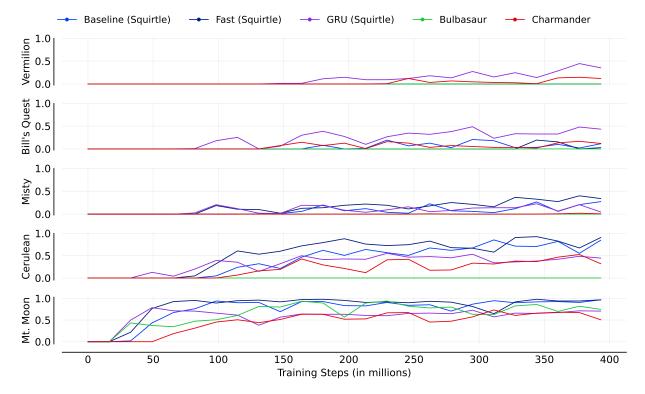


Fig. 3: Mean completion curves for distinct milestones. Beating Misty can be skipped to reach Vermilion City.

Agents may skip Misty since her water-type Pokémon pose a challenge, especially for the fire-type Charmander, which has just a 2% completion rate. Fast (33%), baseline (27%), and GRU (21%) perform better, but the rates are low compared to reaching Cerulean City. While Fast and baseline excel early on, the GRU agent performs better on Bill's quest and on reaching Vermilion City. Bill's quest is not trivial, requiring three basic interactions, in a specified order, without leaving his house. The events reset if the house is exited before completion of the full sequence. Despite the confined space of the house, and an event flags observation that should convey quest progression to agents, completion still proves too complex for most: for example, 71% of Fast experiments visited Bill's house, yet only 19% of all Fast completed his quest. In contrast, 48% of GRU agents completed the quest, down just 1% from the 49% arriving at Cerulean; this suggests an advantage in retaining task-relevant memory.

B. Discussing Ablations and Human Performance

Table II adds the following experiments to those previously described: individual rewards ablations, starting the agent from a state where it can choose its starter Pokémon, and the results from human playthroughs through Cerulean City. Next, we discuss the most salient findings.

1) Distinct Exploration Strategies are Indispensable: Exploring the world of Pokémon is tremendously difficult. To address this, a naive navigation reward is introduced. Without it, the agent fails to achieve any milestones, yet, when this reward is scaled up by a factor of 10, the agent rarely defeats

Brock. Instead, such agents explore almost exclusively, to the detriment of all other aspects of the game. Even battles and menuing are avoided to conserve steps and thus maximize navigation reward.

Since this reward signal is so simplistic, yet navigation is so crucial, it would be interesting to investigate more advanced exploration methods, such as curiosity-driven rewards [9] or Random Network Distillation [10]. Such methods could provide robustness via reduced reward sensitivity, and by potentially allowing the visited-coordinates binary observation to be obviated.

2) Coping with Extremely Long Horizons: Most MDPs have a terminal state, but in Pokémon Red, we must define this ourselves due to the game's open-ended nature. With a fixed episode length, all data-collecting workers reset simultaneously, repeatedly sampling similar training data. Since the game's episodes are significantly longer than the sampled trajectories, there is a high risk of catastrophic forgetting, where learned policies are lost. To mitigate forgetting, we allowed the episode length to grow dynamically based on progress. Because the agent's decisions are stochastic, workers quickly fall out of sync in their resets, with some agents restarting in Pallet Town, while others continue advancing. While this increases sample diversity, training remains unstable: if Figure 3 rendered standard deviations, the performances across different run types would be difficult to distinguish. Similarly, Table II hints on this instability, showing large standard deviations.

Although agents are competitive in defeating Brock, they

fall significantly behind in reaching Cerulean City: humans require about 11,000 steps on average to reach Cerulean City, while an agent requires effectively double that. The Fast agent stood out among agents as the quickest, requiring just 18,500 steps on average. Since Fast greatly speeds the passage of the dialogs that appear upon NPC or object interaction, and when in battle, and clearing such dialogs requires multiple successive, correct inputs, it is likely that action efficiency plays an important role in completion speed. A short dialog takes at least 10 steps; evolution, another dialog, takes about 15 steps. Here, steps can be saved by pressing B, which aborts the evolution. Although human players can use this advantageously, agents only seem to stochastically abort, suggesting an insouciance towards wasted steps. Since addressing such inefficiency would require additional engineering, such as reading memory locations of the game to determine when decision-making is relevant, while our action execution method is precise for overworld navigation, it remains coarse for battles and dialogues. One approach to handle unnecessary actions could be to add a time dimension to the action space. In addition to selecting a button to press, the agent would decide when to act or how long to keep the button pressed. This would allow the agent to learn when to act and when to observe by itself [11].

Finally, the suitability of the training method itself is questionable. Maximizing the discounted cumulative reward encourages minimizing time, introducing bias. This bias is evident in the Choose Starter experiment, where the agent consistently picks Charmander, despite Squirtle being the most effective choice against Brock. The agent favors Charmander because Squirtle and Bulbasaur require one and two extra actions, respectively, to reach in the overworld compared to Charmander. Since picking a starter immediately grants an event and level reward, Charmander is chosen simply because its reward is triggered sooner. By the time the agent reaches Brock, it becomes clear that Charmander is an inferior choice; however, defeating Brock takes over 3,000 steps on average—vastly exceeding both the sampled trajectory length of 2,048 steps and the effective horizon of approximately 333 steps (calculated as $\frac{1}{1-0.997}$). To address this, hierarchical training methods, such as Director [12], are promising. In this context, a sequence of primitive actions (i.e., button presses) can be combined to form high-level actions. Choosing the starter Pokémon could be such an action, which would remove the undesirable time bias. However, learning high-level actions from scratch remains a significant challenge.

3) Heal Reward Exploitation: The heal reward encourages exploitation, especially when Bulbasaur is chosen as the starter. The agent's policy converges to a strategy of battling wild Zubats in Mt. Moon, using Bulbasaur's Leech Seed to restore health while Zubat heals with Leech Life. This results in nearly endless battles, providing frequent heal rewards. Before Bulbasaur faints, the agent visits the Poké Center for another heal reward.

Similar exploitation occurs in the Charmander, GRU, and *Choose Starter* experiments, where at least one out of five runs

involves battling wild Pokémon near Pallet Town and using the player's mother to heal. This strategy leads to suboptimal performance, with runs failing to exceed an 80% success rate against Brock (Table II) and high standard deviations in heals and Poké Center visits (Table III). Removing the heal reward generally reduces performance, as the Poké Center is no longer used as a respawn point. With an average of 8–9 blackouts, the agent requires additional steps to retrace its path from Pallet Town to the point of fainting (e.g., Cerulean City).

Bulbasaur is an exception, as the agent reaches Cerulean City without the heal reward, though it is less effective than Charmander or Squirtle. Bulbasaur may appear advantageous due to its grass-type advantage against Brock, but Bulbasaur learns Vine Whip, its first contact Grass-type attack, comparatively late, at level 15. And, this move has a pitiful 10 PP. In contrast, Squirtle learns Bubble at level 8, and Charmander learns Ember at level 9; both moves have higher PP than Vine Whip, making them more suitable for the lengthy, battle-heavy traversal of Mt. Moon. But such nuance is lost on the agent: the starter Pokémon is chosen very early in the episode, and connecting this choice to distal outcomes would require significant engineering.

4) Reward Shaping Introduces Vulnerabilities: Beyond healing and navigation, event and level rewards also shape the agent's behavior in unintended ways.

Event rewards serve as breadcrumbs, guiding the agent through the storyline. Some, like retrieving a potion from an NPC, are optional, while others, such as earning gym badges, are mandatory. That said, the agent frequently skips Misty's badge in Cerulean City while progressing to Vermilion City, suggesting that critical events may require stronger feedback.

Ablating the level reward yields the best ablation performance (Table II), but introduces side effects as well. The agent's highest-level Pokémon averages 27.45—suprisingly comparable to the baseline, and just two levels below the *Fast* run. To explain this performance, we hypothesize that ablating level reward might allow the agent to 'focus' on other sources of rewards, namely, event rewards. Since most events are trainer battles, which are early-game dense and XP-rich, this agent progresses readily within the scope of this paper by having sufficient level and incentive to readily complete the early game objectives. While all reward sources influence the agent's performance, other rewards have a debatably-positive effect on the agent.

An agent starting with Squirtle (baseline) typically avoids healing exploits, whereas the Charmander agent heals heavily. But, ablate the heal reward, and the agent explores fecklessly, visiting 5637 coordinates on average. Despite varying positive and side effects, no reward nor ablation manages to address the challenge of reaching the third gym in Vermilion City, which is blocked by a cuttable tree (Figure 1f). This presents a significant exploration challenge: the agent must obtain a Pokémon capable of learning Cut, acquire the necessary item, HM01, navigate the UI to use it, and correctly execute the move in the overworld. Until to this point in the game, none of these actions were necessary. Clearly, solving this requires

Mt. Moon, arriving at Cerulean City, defeating Misty, solving Bill's quest, and completing Cerulean City by defeating the Rocket Thief. The best values in each column are highlighted in bold. The Fast experiment saves time (i.e., steps) by setting the game's text speed to fast and disabling battle animations. The Choose Starter experiment does not skip the Oak parcel quest at the beginning of the game, instead starting at the beginning where the player still needs to selects their Rows are sorted in descending order by the number of completed milestones. The tracked milestones include reaching Viridian Forest, defeating Brock, reaching TABLE II: Performance results of various ablation experiments. For each metric, the best data point—determined by milestone completion over time—is selected. first Pokémon.

Experiment	Milestones	Milestones Beat Brock Mt Moon	Mt Moon	Cerulean City Beat Misty	Beat Misty	Bill Quest	Cerulean Done	Beat Brock Steps	Cerulean Done Beat Brock Steps Cerulean City Steps Cerulean Done Steps	Cerulean Done Steps
Human	7.00 ± 0.0	1.00 ± 0.0	1.00 ± 0.0	1.00 ± 0.0	1.00 ± 0.0	1.00 ± 0.0	1.00 ± 0.0	5403 ± 3824	11188 ± 5224	16842 ± 7413
Baseline $-R_{lvl}$	4.60 ± 0.7	0.95 ± 0.0	0.91 ± 0.1	0.79 ± 0.1	0.31 ± 0.3	0.32 ± 0.4	0.03 ± 0.1	6117 ± 1758	21352 ± 3898	37938 ± 0
Fast	4.40 ± 0.8	0.98 ± 0.0	$\boldsymbol{0.98 \pm 0.0}$	0.93 ± 0.1	$\boldsymbol{0.33 \pm 0.4}$	0.19 ± 0.4	0.00 ± 0.0	3753 ± 990	18523 ± 4280	1
Baseline	4.19 ± 0.2	0.99 ± 0.0	0.97 ± 0.0	0.85 ± 0.1	0.27 ± 0.3	0.11 ± 0.2	0.00 ± 0.0	5587 ± 1235	25299 ± 4950	1
Baseline $-R_{heal}$	3.94 ± 0.2	0.99 ± 0.0	0.97 ± 0.0	0.91 ± 0.1	0.07 ± 0.1	0.00 ± 0.0	0.00 ± 0.0	5497 ± 503	28813 ± 2880	I
GRU		0.76 ± 0.4	0.71 ± 0.4	0.49 ± 0.4	0.21 ± 0.3	$\boldsymbol{0.48 \pm 0.4}$	0.21 ± 0.3	5624 ± 780	27118 ± 4236	46532 ± 2249
$Fast-R_{heal}$	3.77 ± 0.5	0.97 ± 0.0	0.96 ± 0.0	0.73 ± 0.3	0.12 ± 0.2	0.00 ± 0.0	0.00 ± 0.0	6531 ± 1185	28871 ± 1362	1
Charmander – Rheal		0.94 ± 0.0	0.92 ± 0.0	0.88 ± 0.1	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	11694 ± 1439	30369 ± 4828	1
Choose Starter		0.79 ± 0.4	0.79 ± 0.4	0.75 ± 0.4	0.29 ± 0.4	0.00 ± 0.0	0.00 ± 0.0	6485 ± 578	20279 ± 579	I
Choose Starter $-R_{heal}$		0.97 ± 0.0	0.96 ± 0.0	0.45 ± 0.3	0.01 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	13871 ± 2991	33565 ± 13303	1
Bulbasaur $-R_{heal}$		0.97 ± 0.0	0.95 ± 0.0	0.26 ± 0.2	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	13433 ± 950	45943 ± 3569	1
Charmander		0.69 ± 0.4	0.67 ± 0.3	0.53 ± 0.4	0.02 ± 0.0	0.17 ± 0.3	0.02 ± 0.0	7029 ± 1132	21309 ± 5915	48544 ± 0
Bulbasaur	2.89 ± 0.1	0.95 ± 0.0	0.94 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	9929 ± 841	1	1
Baseline $+ 10R_{nav}$	1.04 ± 0.1	0.04 ± 0.1	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	11559 ± 0	ı	ı
Baseline $-R_{nav}$	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	1	ı	I

TABLE III: Policy metrics describing the agent's behavior. Rows are ordered as in Table II. The best value in each column is highlighted in bold. The species entropy refers to the probability distribution of catching distinct species. As the Fast agent usually does not catch Pokémon, but may evolve Squirtle to Wartortle, the entropy is lowest for these runs. The human playthroughs are more diverse in their party setup as indicated by the largest entropy.

Experiment	Events Completed	Events Completed Poké Center Visits Num. Heals	Num. Heals	Black Outs	Max Party Level	Visited Coordinates	Party Size	Species Entropy
Human	50.47 ± 9.1	12.53 ± 7.7	22.97 ± 11.3	2.23 ± 2.1	26.20 ± 2.1	1997 ± 568	2.07 ± 1.2	3.53 ± 0.0
Baseline $-R_{lvl}$	57.73 ± 11.4	20.29 ± 24.2	11.73 ± 4.0	9.01 ± 3.2	27.45 ± 2.8	4560 ± 1015	1.58 ± 0.5	1.24 ± 0.2
Fast	57.21 ± 5.4	111.46 ± 133.9	47.98 ± 66.6	9.68 ± 1.7	29.43 ± 1.9	3392 ± 744	3.49 ± 0.2	2.54 ± 0.1
Baseline	51.78 ± 2.1	19.09 ± 19.6	12.81 ± 5.9	9.08 ± 2.5	27.52 ± 3.2	4275 ± 1406	3.48 ± 0.3	2.61 ± 0.1
Baseline $-R_{heal}$	43.21 ± 2.0	0.04 ± 0.1	1.45 ± 0.7	9.74 ± 1.2	23.38 ± 1.1	5637 ± 368	3.72 ± 0.2	2.27 ± 0.2
GRU	47.20 ± 21.3	18.16 ± 25.3	23.09 ± 23.3	8.25 ± 5.2	22.90 ± 9.4	2534 ± 1421	2.72 ± 0.9	1.87 ± 1.0
$Fast-R_{heal}$	41.79 ± 5.3	0.01 ± 0.0	2.70 ± 1.1	9.27 ± 1.8	23.73 ± 0.8	5135 ± 843	3.20 ± 0.3	2.16 ± 0.2
Charmander $-R_{heal}$	45.29 ± 7.2	0.01 ± 0.0	2.68 ± 1.9	8.54 ± 1.6	24.61 ± 2.5	5399 ± 511	3.32 ± 0.3	2.09 ± 0.2
Choose Starter	48.42 ± 17.7	44.32 ± 36.1	33.75 ± 34.1	9.14 ± 5.7	24.72 ± 8.7	2577 ± 1160	1.79 ± 0.4	1.25 ± 0.6
Choose Starter $-R_{heal}$	41.49 ± 3.0	0.01 ± 0.0	5.75 ± 2.4	9.73 ± 2.2	24.88 ± 2.2	4527 ± 371	1.86 ± 0.8	1.30 ± 0.6
Bulbasaur — R_{heal}	36.55 ± 2.6	0.01 ± 0.0	8.74 ± 4.6	10.00 ± 0.7	21.80 ± 0.8	4787 ± 406	3.43 ± 0.6	2.34 ± 0.2
Charmander	46.52 ± 19.7	57.83 ± 57.8	25.87 ± 19.0	6.01 ± 3.5	22.71 ± 8.6	2871 ± 1415	2.95 ± 1.0	2.22 ± 1.1
Bulbasaur	30.32 ± 2.2	88.15 ± 20.7	399.33 ± 177.1	7.76 ± 3.5	20.63 ± 3.8	2283 ± 279	3.17 ± 0.5	2.37 ± 0.3
Baseline $+ 10R_{nav}$	16.39 ± 1.2	0.00 ± 0.0	0.02 ± 0.0	3.49 ± 2.5	9.05 ± 1.7	3284 ± 364	1.03 ± 0.0	0.18 ± 0.2
Baseline - R _{nav}	13.40 ± 1.2	0.00 ± 0.0	0.00 ± 0.0	0.00 ± 0.0	6.00 ± 0.0	35 ± 15	1.00 ± 0.0	0.00 ± 0.0

an additional reward signal to incentivize all the behaviors associated with Cut. Indeed, many such challenges exist in Pokémon Red, and further refining the reward function risks introducing other vulnerabilities to be exploited and side effects to be realized, highlighting the environment's complexity.

5) Limitations: While Pokémon Red serves as a valuable research environment, it comes with several limitations. First, the game is closed-source and requires owning a legal copy. Once a digital copy is available for training, additional engineering effort is needed to extract relevant information from the Game Boy emulator's RAM.

Second, the prolonged episode length significantly increases evaluation time. Running evaluation episodes takes as much wall-time as the training runs themselves, with a single training run lasting approximately 36 hours. This extended horizon also presents challenges for agents utilizing Recurrent Neural Networks. Processing long observation sequences—2048 steps in our case—leads to substantial VRAM overhead. As such, in lieu of additional optimization, the GRU experiments were run on CPU only, with an average runtime of 24 days per training run.

Furthermore, the vast number of possible gameplay strategies contributes to the curse of dimensionality. Comparing agents across numerous factors makes comprehensive evaluation labor-intensive. For example, analyzing item usage and move selection falls outside the scope of this paper, as does manually reviewing thousands of episodes.

Addressing these limitations may require refining the toolset, which has been effective for other environments.

V. RELATED WORK

To our knowledge, no peer-reviewed publications have explored training agents on Pokémon Red or its successors. Recent works have, however, introduced Large Language Model (LLM) agents. A preprint presented an agent focused solely on Pokémon battles, achieving human-parity performance [13]. The startup nunu.ai demonstrated an agent that reached the third badge in Pokémon Emerald, relying on human intervention and a pathfinding algorithm. More recently, Anthropic announced efforts to develop an agent capable of earning the third badge in Pokémon Red [14]. This is particularly noteworthy, as Gym 3 requires both using Cut and solving a hidden button puzzle. Although details on these agents remain limited, the increasing interest in Pokémon as an LLM testbed highlights its growing relevance in AI research. A promising direction for future work is an agent similar to Minecraft Voyager [15], which generates and evaluates code to perform actions.

VI. CONCLUSION AND OUTLOOK

Our simple DRL baseline trains agents to complete an initial portion of Pokémon Red, serving as stepping stone for tackling its extensive challenges. Key tasks like cutting trees and solving puzzles remain unexplored by our agents but pose significant hurdles, potentially requiring engineered solutions such as reward shaping or curriculum learning, thus further

diverging from human perception and play. We did not address generalization, as agents started from a fixed state while influencing RNG through their actions. ROM randomizers could improve robustness by varying start conditions, such as the choice of starter Pokémon. For handling long horizons, hierarchical DRL may reduce unnecessary actions, while advanced exploration methods could replace naive navigation rewards and binary image observations. Overall, we anticipate Pokémon Red will play an essential role in future research.

ACKNOWLEDGMENT

This research is supported by computing time from the Paderborn Center for Parallel Computing and the Linux-HPC-Cluster at TU Dortmund. We also thank Keelan Donovan and Sky (Discord username) for their discussions and reviews, PufferLib [16] for RL expertise, and Ryan Sullivan and Joseph Suarez for manuscript review.

REFERENCES

- C. Berner, G. Brockman, B. Chan, V. Cheung, P. Debiak, C. Dennison,
 D. Farhi, Q. Fischer, S. Hashme, C. Hesse, R. Józefowicz, S. Gray,
 C. Olsson, J. Pachocki, M. Petrov, H. P. de Oliveira Pinto, J. Raiman,
 T. Salimans, J. Schlatter, J. Schneider, S. Sidor, I. Sutskever, J. Tang,
 F. Wolski, and S. Zhang, "Dota 2 with large scale deep reinforcement learning," CoRR, vol. abs/1912.06680, 2019.
- [2] H. Küttler, N. Nardelli, A. H. Miller, R. Raileanu, M. Selvatici, E. Grefenstette, and T. Rocktäschel, "The NetHack Learning Environment," in NIPS, 2020.
- [3] L. Fan, G. Wang, Y. Jiang, A. Mandlekar, Y. Yang, H. Zhu, A. Tang, D. Huang, Y. Zhu, and A. Anandkumar, "Minedojo: Building openended embodied agents with internet-scale knowledge," in NIPS, 2022.
- [4] P. Whidden, "Training ai to play pokemon with reinforcement learning," 2023, accessed: 2025-02-25. [Online]. Available: https://www.youtube.com/watch?v=DcYLT37ImBY
- [5] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *CoRR*, vol. abs/1707.06347, 2017.
- [6] R. S. Sutton and A. G. Barto, Reinforcement Learning: An Introduction, 2nd ed. MIT Press, 2018.
- [7] M. Pleines, M. Pallasch, F. Zimmer, and M. Preuss, "Memory gym: Partially observable challenges to memory-based agents," in *ICLR*, 2023.
- [8] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A. Riedmiller, A. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, pp. 529–533, 2015.
- [9] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, "Curiosity-driven exploration by self-supervised prediction," in *ICML*, 2017.
- [10] Y. Burda, H. Edwards, A. J. Storkey, and O. Klimov, "Exploration by random network distillation," in *ICLR*, 2019.
- [11] H. Zhou, A. Huang, K. Azizzadenesheli, D. Childers, and Z. Lipton, "Timing as an action: Learning when to observe and act," in AISTATS, 2024.
- [12] D. Hafner, K.-H. Lee, I. Fischer, and P. Abbeel, "Deep hierarchical planning from pixels," in NIPS, 2022.
- [13] S. Hu, T. Huang, and L. Liu, "Pokellmon: A human-parity agent for pokemon battles with large language models," CoRR, vol. abs/2402.01118, 2024.
- [14] Anthropic, "Claude's extended thinking," 2025, accessed: 2025-02-25. [Online]. Available: https://www.anthropic.com/news/ visible-extended-thinking
- [15] G. Wang, Y. Xie, Y. Jiang, A. Mandlekar, C. Xiao, Y. Zhu, L. Fan, and A. Anandkumar, "Voyager: An open-ended embodied agent with large language models," *CoRR*, vol. abs/2305.16291, 2023.
- [16] J. Suarez, "Pufferlib: Making reinforcement learning libraries and environments play nice," 2024, accessed on March 12, 2025. [Online]. Available: https://arxiv.org/abs/2406.12905