# DEEP REINFORCEMENT LEARNING FOR POKEMON BATTLING
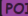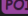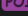
**Kevin Zhang**

PURDUE UNIVERSITY® | Department of Computer Science
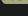
## Pokemon Battles offer unique challenges

- Stochastic: Accuracy, Status, Crits

- Partially-Observable: Stats, abilities, moves

- Very large state/search space: Embedding of above

- Additionally- Information may be entirely known by "testing" during battle

# *Pokemon Showdown Environment*

## Turning Pokemon into an OpenAI Gym environment

- Pokemon Showdown offers an open-source method for generating and conducting battles

- Poke-Environment acts as an interface to Python for this data

- Select an action at each state, and use each turn as a transition
  - Follow Gymnasium's specification
  - Integrate Agents trained using Pytorch

# Training Design

**There are many options for training methodology**

- ▪ Algorithm-
  - • Optimization Method: DQN, PPO, GIGAθ, WPL, etc.
  - • Opponent Design: Baseline, Selfplay, etc.

- ▪ Information Structure-
  - • Embedding huge information space

- ▪ Reward Design-

**Algorithm 1:** Pseudo-code for a worker thread running GIGAθ using -greedy exploration.

**Input:** Globally, target network update period $\tau$, on-line and target Q-networks, on-line and target policy networks, on-line and target average policy network weights, exploration rate $\epsilon$, and maximum iterations $T_{max}$. Locally, on-line Q, policy and average policy networks.

1: **for** iteration $T \leftarrow 0, T_{max}$ **do**
2:   Synchronize local on-line networks as copies of global on-line networks
3:   Sample initial state
4:   **repeat**
5:     Execute random action with probability $\epsilon$, otherwise execute action from policy network
6:     Sample new state and reward
7:     Compute Q-network targets with target Q-network
8:     Compute policy networks' targets with GIGA-WoLF's update equations
9:     Compute loss of local on-line Q, policy, and average policy networks
10:    Accumulate gradients by minimizing the loss
11:  **until** terminal state
12:  Update global on-line networks weights with the accumulated gradients of local on-line networks
13:  Synchronize global target networks as copies of global on-line networks every $\tau$ time-steps
14: **end for**

**Output:** A converged Q-network to approximate the value function as $Q(s, a, \theta)$, and a converged policy network to approximate the policy function as $\pi(s, a, \theta_\pi)$.

**Competitive Deep Reinforcement Learning over a Pokémon Battling Simulator** (Simoes et. al)

# Deep Reinforcement Learning

## Q Networks

- Reinforcement Learning is based on the Bellman Equations

$$v_\pi(s) = \mathbb{E}_\pi \left[ R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s \right]$$

$$q_\pi(s, a) = \mathbb{E}_\pi \left[ R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a \right]$$

- Storing these mappings is difficult: Learn them with a neural network
  - State -> Q-values for each action

- Deep Q Network (DQN) – Train a NN to select actions via argmax
  - Improvements such as Replay and Target networks improve stability

    **Human Level Control through Deep Reinforcement Learning**
    (Mnih et. al)

# Strategies

**After integrating with environment, test:**

- Algorithms:
  - DQN, and improvements (n-step learning, priority replay)
  - GIGAθ, WPL

- Embeddings and Reward:
  - Naïve: Only take into account power & hp
  - Hand Designed: Take into account stats of bench and opponent
  - Future- Reward Network

- Opponent
  - Random & MaxPower Baseline
  - Selfplay

# *Results*

## Initial Results-

- DQN with Random and MaxPower Baseline: 5 runs
  - Training over 10,000 env steps

| Training against Max | | | Training against Rand | | |
|---|---|---|---|---|---|
| Agent-Rand | Agent-Max | Agent-Heuristic | Agent-Rand | Agent-Max | Agent-Heuristic |
| 0.9 | 0.38 | 0.1 | 0.58 | 0.2 | 0.02 |
| 0.94 | 0.4 | 0.12 | 0.82 | 0.22 | 0.04 |
| 0.84 | 0.44 | 0.14 | 0.76 | 0.38 | 0.06 |
| 0.9 | 0.56 | 0.12 | 0.64 | 0.26 | 0.06 |
| 0.9 | 0.3 | 0.2 | 0.7 | 0.22 | 0.06 |
| 0.896 | 0.416 | 0.136 | 0.7 | 0.256 | 0.048 |
| 0.035777088 | 0.095289034 | 0.038470768 | 0.09486833 | 0.072663608 | 0.017888544 |

- Beating similar student project's performance
  ~65% - 70% : Kalose et. al.
  https://web.stanford.edu/class/aa228/reports/2018/final151.pdf

(a) Results when trained against a MaxPowerPlayer baseline.
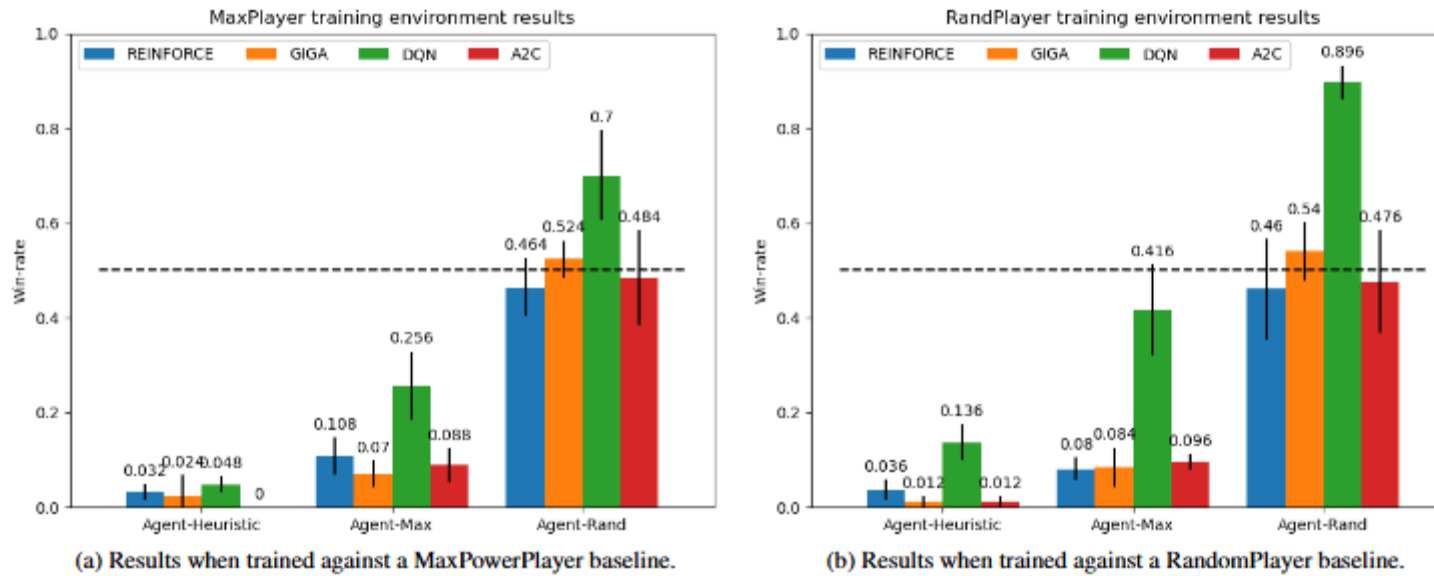
(b) Results when trained against a RandomPlayer baseline.

Figure 5. Final win-rate results for a sample size of n=5 runs. Error bars show sample standard deviation. The x-axis corresponds with the opponent in testing, the two graphs shows the results of models trained against different baseline models. The dashed line shows a win-rate of 0.5 for approximate parity in strength.

| Experiment | Win-rate | | |
| --- | --- | --- | --- |
| | Agent-Rand | Agent-Max | Agent-Heuristic |
| MaxPower DQN | $0.896 \pm 0.036$ | $0.416 \pm 0.095$ | $0.136 \pm 0.0385$ |
| Heuristic DQN | $0.66 \pm 0.105$ | $0.276 \pm 0.056$ | $0.072 \pm 0.023$ |
| Adam DQN | $0.712 \pm 0.107$ | $0.344 \pm 0.0385$ | $0.152 \pm 0.048$ |

Figure 6. Final win-rate results for a sample size of n=5 runs. Adam DQN is trained with an Adam optimizer against the MaxPower baseline agent. Errors show sample standard deviation.

# *Future Work*

## Directions for this Project

- Simple:
  - Develop Embeddings & Reward
  - Implement other algorithms/improvements (Rainbow is all you need)
  - Hyperparameter tuning
- Advanced:
  - Reward network
  - Develop novel training algorithms
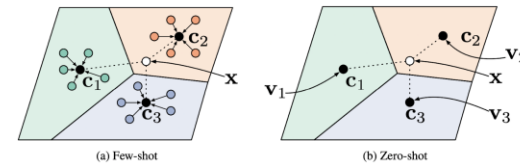    - Potentially tackle IL/BC or Few-shot

Figure 1: Prototypical Networks in the few-shot and zero-shot scenarios. **Left:** Few-shot prototypes $c_k$ are computed as the mean of embedded support examples for each class. **Right:** Zero-shot prototypes $c_k$ are produced by embedding class meta-data $v_k$. In either case, embedded query points are classified via a softmax over distances to class prototypes: $p_\phi(y = k|\mathbf{x}) \propto \exp(-d(f_\phi(\mathbf{x}), \mathbf{c}_k))$.

(a) Few-shot     (b) Zero-shot