

Understanding biological active sensing behaviors by interpreting learned artificial agent policies

Sonja Johnson-Yu¹, Satpreet H. Singh¹, Federico Pedraja², Denis Turcu², Pratyusha Sharma³, Naomi Saphra¹, Nathaniel B. Sawtell², Kanaka Rajan¹
Harvard University¹, Columbia University², MIT³

Abstract

Weakly electric fish, such as *Gnathonemus petersii*, generate pulsatile electric organ discharges (EODs) that enable them to sense their environment through active electrolocation. This plays a crucial role in several key behaviors, such as navigation, foraging, and avoiding predators. While the anatomical and physiological organization of the active electrosensory system has been extensively studied, the contribution of active electrolocation to adaptive behavior in naturalistic settings remains relatively underexplored. Here we present a preliminary *in silico* model of active sensing in electric fish, using a neural network-based artificial agent trained by deep reinforcement learning to perform an analogous active sensing task in a 2D environment. The trained agent recapitulates key features of natural EOD statistics, shows emergent behavioral modularity, and provides intuitions about the representation of key latent variables governing agent behavior, such as energy levels (satiety).

1 Introduction

Weakly electric fishes like *Gnathonemus petersii* use electric pulses, or electric organ discharges (EODs), to actively sense their environment, communicate with each other, and sense their environment based on the EODs of nearby fish (Von der Emde, 1999; Sawtell et al., 2005; Pedraja & Sawtell, 2024). The role that active electrolocation plays in the goal-oriented behaviors of fish is less well understood compared to our extensive knowledge of the physiology of the neural mechanisms responsible for EOD generation. This knowledge gap is due to the difficulty of designing naturalistic yet well-controlled studies that capture the complexity of the animals' sensory ecology and behavioral repertoire.

In recent years, neural network-based artificial agents trained to perform different tasks have emerged as powerful tools to model animal behaviors and neural computations (Haesemeyer et al., 2019; Singh et al., 2023). By transforming sensory inputs into motor outputs similar to those of real animals, such models offer insight into the neural and cognitive processes underlying animal behaviors. They also enable flexible *in silico* experimentation while being fully observable, allowing hypothesis testing where experimental data collection is challenging.

Here, we present preliminary results from a biologically-constrained artificial agent trained by deep reinforcement learning (DRL) to perform an active-sensing foraging task in a 2D environment, analogous to weakly electric fish behavior.

2 Environment and Agent

2.1 Overview

Inspired by lab experiments on *Gnathonemus petersii*, we train our agents in simulated 2D tanks of size 60 cm x 60 cm (Fig 1a). Simulations are initialized with n food items placed uniformly at

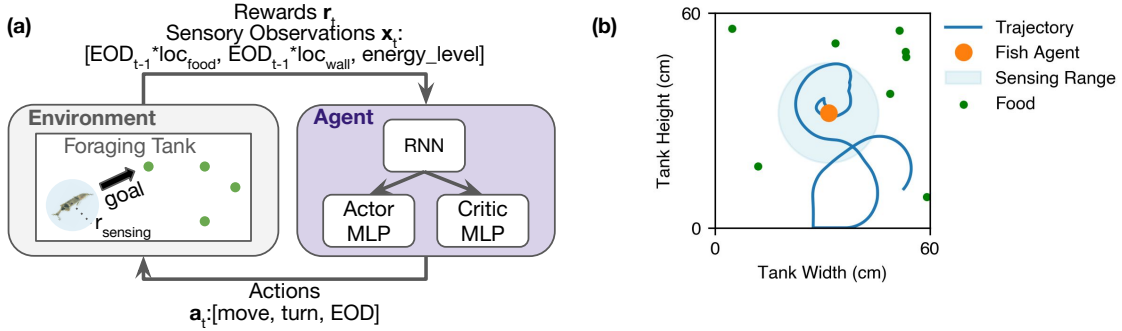


Figure 1: **(a) The agent works in tandem with the tank simulation environment to learn an efficient foraging policy** (billycorgan84, 2009). At each timestep, the agent receives sensory observations and rewards from the environment and then selects its actions. If the agent emitted an EOD in the previous timestep, it can observe the location of the nearest food and wall within its sensing range. The agent uses a recurrent neural network (RNN) to infer the environmental (‘belief’) state, selects an action using an Actor multilayer perceptron (MLP), and estimates the action’s value with a Critic MLP. **(b) Example trajectory from a trained agent.** Food is distributed uniformly at random throughout the 60cm \times 60cm tank. The agent can sense food within its sensing range (radius=14cm).

random. The position and orientation of a single agent are also initialized uniformly at random. At each timestep, the agent observes the egocentric vector distance of the nearest food item and the nearest wall within its sensing range. It can also observe its internal energy levels ($e \in [0, 1]$), which increase every time it eats food, but otherwise decrease linearly with time and activity levels. At each timestep, it decides how much to move forward, how much to turn, and whether or not to emit an EOD. The agent is rewarded for eating food, penalized for both starvation and overeating, and has a baseline metabolic cost associated with staying alive. The agent (Fig 1a, right) consists of a recurrent neural network (RNN) (Rajan et al., 2016) followed by parallel two-layer Actor and Critic Multi-Layer Perceptrons (MLPs). The former selects the agent’s actions, and the latter estimates the value of actions during training using policy gradients (Ni et al., 2021). All layers are 64-units wide, with \tanh nonlinearities. We constrain the agent’s maximum linear and angular velocities and accelerations to match experimental data collected from an electric fish in an identically-sized tank. Simulations are run at ≈ 83 FPS to enable a minimum SPI of 12ms, as is observed in lab experiments. For simplicity, here, the agent actions and observations are deterministic.

2.2 Environment

The environment is modeled as a Partially Observable Markov Decision Process (POMDP), a framework that models scenarios where the agent has incomplete information about the true state of the environment. The POMDP is specified as a collection of possible states \mathcal{S} , actions \mathcal{A} , observations Ω , transition probabilities \mathcal{T} , observation probabilities \mathcal{O} , and rewards $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. In the tank environment, the state s_t is a vector of the agent’s position and orientation in continuous x, y space, the agent’s linear and angular velocities, the agent’s energy level $e \in [0, 1]$, and the locations of all uneaten food items. At each time step, the agent chooses an action a_t composed of linear acceleration, angular acceleration, and whether or not to emit an electric pulse. While the agent is allowed to observe its energy level at every time step, the other observations (Fig 1a, top) are conditional on several factors. If the agent has emitted an electric pulse in the previous time step, then it has the opportunity to observe the location of the nearest food item and the nearest wall, if these fall within the agent’s sensing radius (14 cm). The locations are observed as an egocentric angle and a normalized proximity (proximity = $1 - \frac{\text{distance}}{\text{sensing radius}}$) for both the nearest food and the nearest point on the wall. If the items do not fall within the sensing radius, or if the agent did not

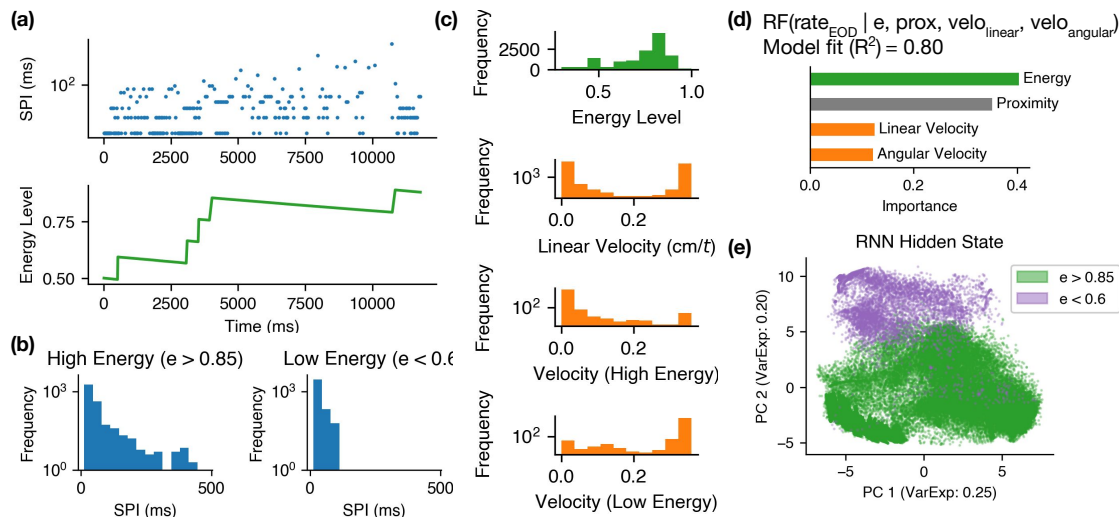


Figure 2: **(a) Example sequential pulse intervals (SPIs) and energy of a trained agent over a 1500-timestep episode.** A sequential pulse interval is the length of time between two EODs. Periods of repeated low SPIs (frequent EODs) correlate with vigorous foraging behavior, seen in the step-wise increases in energy between 0-5000ms. High energy (satiated) behavior correlates with higher SPIs (infrequent EODs) after 5000ms. **(b) Distribution of SPIs when the agent has high energy vs. low energy**, from 30 episodes of 1500 timesteps each. When the agent’s energy is high, its discharge patterns show a wide range of SPIs, including high SPIs (infrequent EODs). When the agent’s energy (satiety) is low, its SPIs are low because it is actively foraging with more frequent EOD discharges. Each SPI incurs a metabolic cost, so it is notable that the low-energy agent pulses frequently. This indicates that the low-energy agent prioritizes finding food to avoid starvation, rather than conserving energy. **(c) Distribution of agent energy levels and linear velocities** (in cm/timestep) across 30 episodes of 1500 timesteps each. The agent tends to maintain its energy at a “set point” close to full. Above this set point, the agent is penalized for overeating. The agent’s linear velocity is bimodal (not swimming vs. swimming vigorously). High energy levels (high satiety) correlate with low velocity, and conversely, low energy levels correlate with high velocity. The agent’s energy level appears to influence its locomotion strategy. When the agent’s energy level (satiety) is high, it does not need to eat more food and swims slowly. High velocity often corresponds to an agent motivated to eat more food and gain energy. **(d) Feature importance from a 100-tree random forest predicting EOD rate.** Agent energy, followed by proximity to food, is the most important predictor of EOD rate. **(e) Principal component analysis (PCA) of the RNN’s hidden states**, from 30 episodes of 3000 timesteps each. The hidden state output by the RNN at each timestep can be interpreted as a low-dimensional “summary” of the agent’s belief about the state of the environment. We observe a transition along the 2nd principal component between the hidden states corresponding to a low-energy agent vs. those of a high-energy agent. This indicates that energy may be a latent variable that plays an important role in determining the agent’s actions.

emit a pulse in the previous timestep, the agent observes a vector of zeros. The transitions of the environment are deterministic and the observation probabilities are 1. Lastly, the agent’s goal is eat a sufficient amount of food and thereby avoid starvation. The agent’s reward for each timestep is specified as such:

$$r_t = r_{\text{food eaten}} - r_{\text{metabolism}}(a_t) - r_{\text{overeate}} - r_{\text{starvation}}$$

where $r_{\text{metabolism}}(\text{accel}_{lin}, \text{accel}_{ang}, \text{pulse}) = a \cdot \text{accel}_{lin} + b \cdot \text{accel}_{ang} + c \cdot \text{pulse} + d$, and a, b, c, d are tunable hyperparameters.

This reward function incentivizes the agent to forage for food items efficiently without “overeating,” or going above the allowed energy level.

The simulated tank environment, pictured in Fig 1b, built on OpenAI’s Gym module, is 60 cm \times 60 cm and is initialized with $n = 40$ food items located uniformly at random throughout the tank.

2.3 Agent

The agent learns via PPO, gradient-based algorithm to learn a policy, which can be challenging when state, action, observation, and reward spaces are all continuous. At each timestep t , the agent receives sensory observations x_t and rewards r_t from the environment. The agent’s architecture (Fig 1a, right) is composed of a recurrent neural network (RNN) (Rajan et al., 2016) with a 64-unit hidden layer, which is used to infer the state from the observation x_t . The RNN’s hidden state is then passed both to an Actor Multi-Layer Perceptron (MLP), which outputs an action a_t , and to a Critic MLP, which outputs an estimate of the expected returns from the given state v_t . During the process of training, v_t is used to update the PPO algorithm by comparing the predicted and actual returns.

2.3.1 Biological Constraints

We impose biological constraints on the agent and environment in order to make the agent’s learned policy more realistic. For example, each timestep is equivalent to 12 ms, the minimum latency for consecutive pulses in *Gnathonemus petersii*. Additionally, the agent’s maximum linear and angular velocities and accelerations were learned from the trajectory data of a single elephantfish in a 60 cm \times 60 cm tank.

3 Results

Trained agents successfully electrolocate food items while producing movement trajectories (Fig. 1b) and EOD transcripts (Fig. 2a) that resemble experimental data from real fish. Two behavior modes (Fig 2b), namely “resting” and “active foraging”, are also observed, similar to those observed in real fish (von der Emde, 1992). We also observe that the trained agent learns a “homeostatic drive” to maintain its energy levels slightly below the maximum possible (Fig. 2c). Additionally, we find that the high energy (satiated) state is correlated with low linear velocity and vice versa (Fig. 2c). Furthermore, the low- and high- energy modes are observable in the RNN’s hidden state activities (Fig. 2e).

4 Conclusions

In summary, our *in silico* artificial neural-network agent model recapitulated key features of active sensing behavior in electric fish, including similar EOD statistics and emergent ‘active’ and ‘rest’ behavioral states. Our preliminary analysis of the neural activity underlying the learned policy revealed a potential latent variable, i.e. energy level (satiety), that seems to govern agent behavior. In the future, we hope to cross-pollinate interpretability techniques being developed in Computer Science and Computational Neuroscience to further analyze agent behavior and neural-activity. We also plan to explore the role of EODs in more complex tasks involving cooperation and competition between multiple identical and diverse agents.

5 Others

Broader Impact Statement

Our research contributes to the acceleration of hypothesis generation in neuroscience by leveraging *in-silico* experimentation, paving the way for a deeper understanding of neural processes in biological systems. Furthermore, techniques developed in Computational Neuroscience might potentially inspire new methods for agent interpretability in Computer Science.

References

- billycorgan84. Gnathonemus_petersii, 2009. URL https://commons.wikimedia.org/wiki/File:Gnathonemus_petersii.jpg.
- Martin Haesemeyer, Alexander F Schier, and Florian Engert. Convergent temperature representations in artificial and biological neural networks. *Neuron*, 103(6):1123–1134, 2019.
- Tianwei Ni, Benjamin Eysenbach, and Ruslan Salakhutdinov. Recurrent model-free rl is a strong baseline for many POMDPs. *arXiv preprint arXiv:2110.05038*, 2021.
- Federico Pedraja and Nathaniel B Sawtell. Collective sensing in electric fish. *Nature*, pp. 1–6, 2024.
- Kanaka Rajan, Christopher D Harvey, and David W Tank. Recurrent network models of sequence generation and memory. *Neuron*, 90(1):128–142, 2016.
- Nathaniel B Sawtell, Alan Williams, and Curtis C Bell. From sparks to spikes: information processing in the electrosensory systems of fish. *Current opinion in neurobiology*, 15(4):437–443, 2005.
- Satpreet H Singh, Floris van Breugel, Rajesh PN Rao, and Bingni W Brunton. Emergent behaviour and neural dynamics in artificial agents tracking odour plumes. *Nature Machine Intelligence*, 5(1):58–70, 2023.
- Gerhard von der Emde. Electrolocation of capacitive objects in four species of pulse-type weakly electric fish: II. electric signalling behaviour. *Ethology*, 92(3):177–192, 1992.
- Gerhard Von der Emde. Active electrolocation of objects in weakly electric fish. *Journal of experimental biology*, 202(10):1205–1215, 1999.