

# Fido: A Universal Robot Control System Using Reinforcement Learning with Limited Feedback

Joshua Gruenstein   Michael Truell

## Control System Objectives

Fido was created to fulfill the following goals:

- **Trainability:** Allow both human and autonomous training and retraining rather than programming
- **Universality:** Run on any robot and perform any task, even without prior knowledge of the host
- **Speed:** Require few learning iterations and low latency for quick and efficient training

To achieve this we developed a novel machine learning algorithm utilizing wire-fitted Q-Learning with a neural network, a probabilistic action selection model, and history sampling. The control system could be trained at a rate **3 times greater** than the industry standard, allowing practicality over traditional pre-programmed control systems.

## System Overview

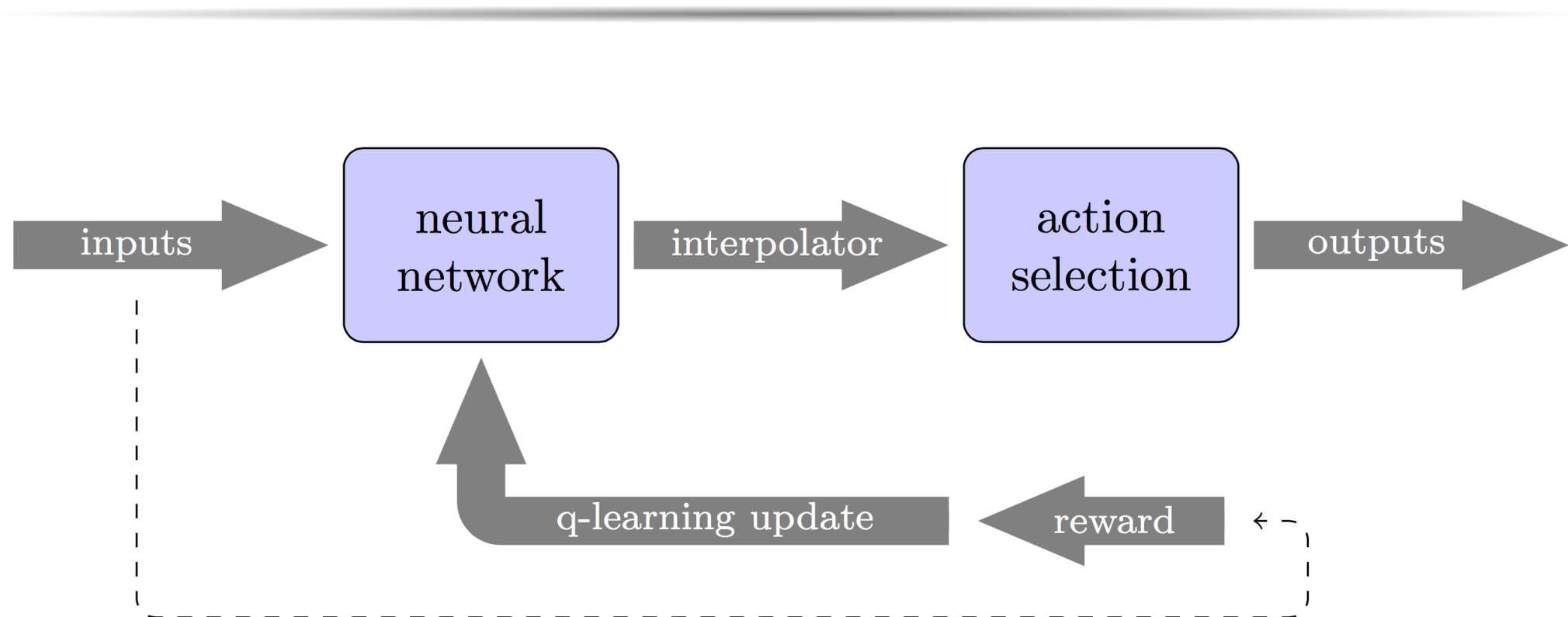


Figure 1: Control System Diagram

- From a macro perspective, Fido can be viewed as a “black box:” inputs go in and outputs go out, Fido must optimize the relationship of inputs to outputs to maximize reward
- **Reward system:** Trying to determine the expected reward for an action in a given state based on past reward received
  - Must have a scalable, performance-optimized way of storing past state-reward sets and detecting patterns

## Reinforcement Learning

- **Q-Learning:** Develops a function that intakes a state-action pair and outputs expected utility
- Ordinarily, the Q-function is modeled by storing state-action pairs in a table
  - Is impractical for large state spaces, use a function approximator instead: **Artificial Neural Networks**
- Usually discrete: no relation made between states or actions
  - Can be optimized by coupling a wire-fitted interpolator with our neural network (**Wire-Fitted Q-Learning**)
- Cannot just pick the action with the greatest expected utility: must “explore” to be trainable and re-trainable
  - Use a **Boltzmann Probability Distribution** selection policy

## Action Selection

- Function approximators modeled after nature with the capability to take in a large number of inputs, parallelly process them, and produce a set of outputs

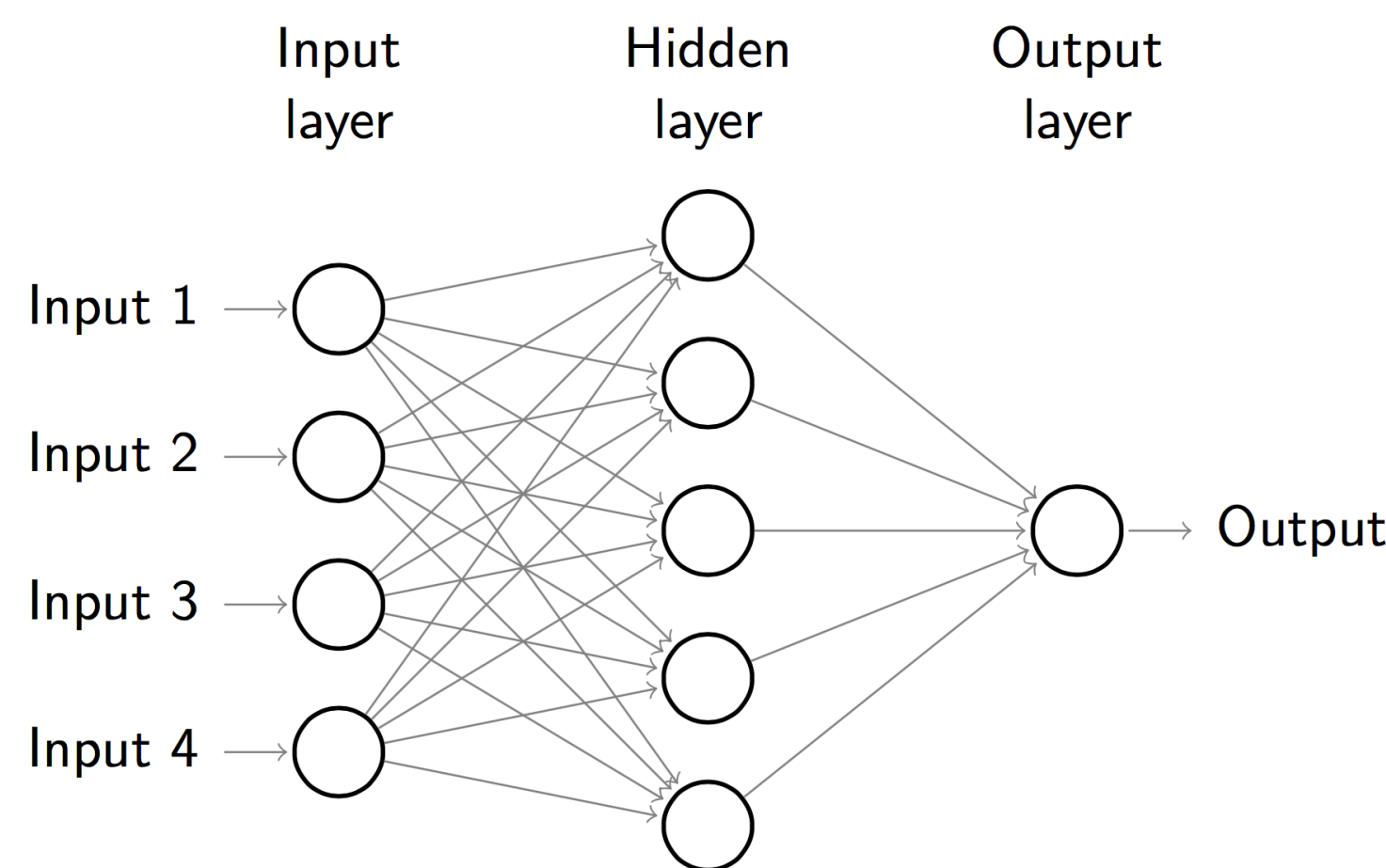


Figure 2: Single Output Feed-forward Neural Network

## Training

- **Neural Networks:** Function approximators modeled after nature with the capability to take in a large number of inputs, parallelly process them, and produce a set of outputs

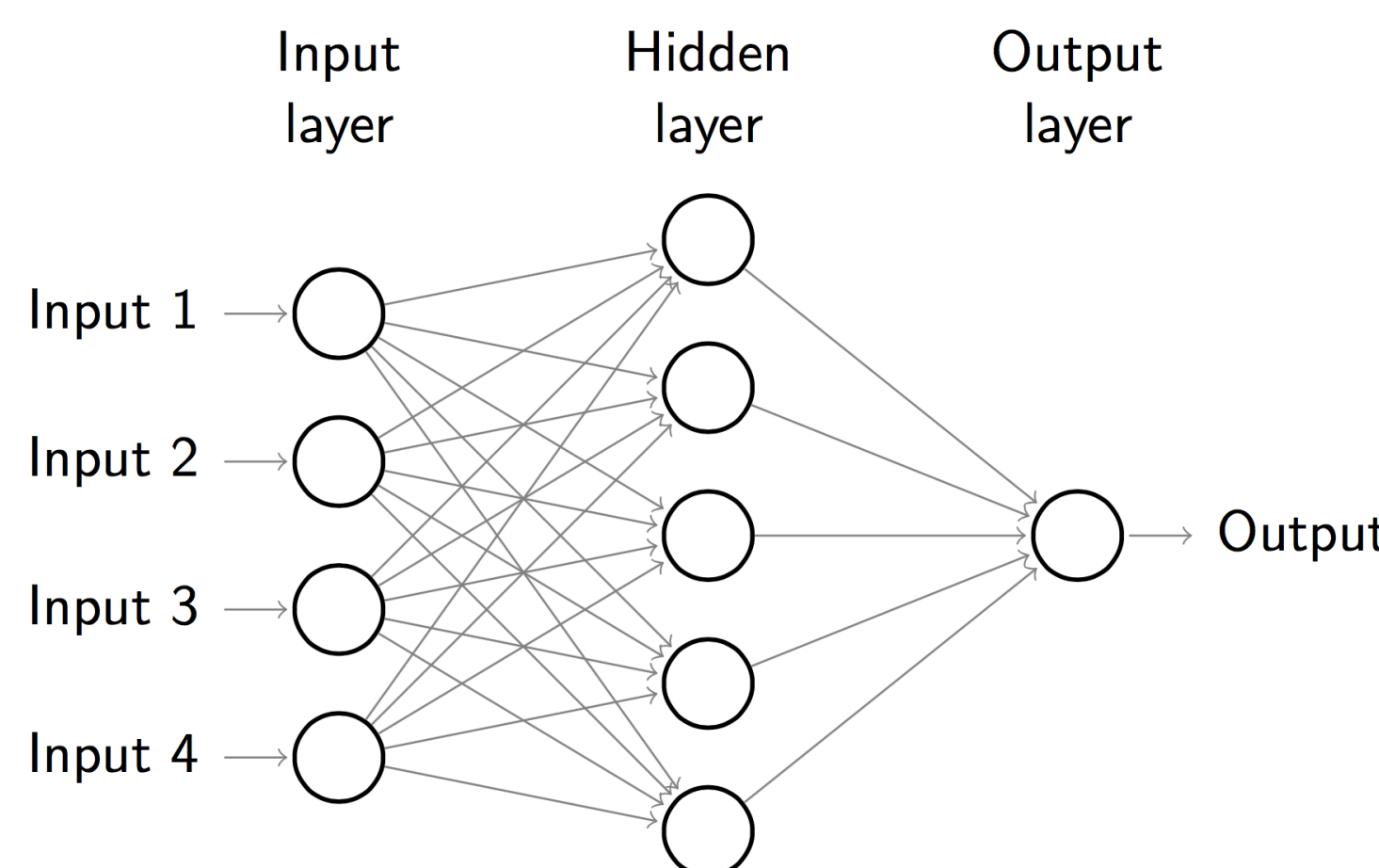


Figure 3: Single Output Feed-forward Neural Network

- **History Sampling:** Develops a function that intakes a state-action pair and outputs expected utility
- Ordinarily, the Q-function is modeled by storing state-action pairs in a table
  - Is impractical for large state spaces, use a function approximator instead: **Artificial Neural Networks**
- Usually discrete: no relation made between states or actions
  - Can be optimized by coupling a wire-fitted interpolator with our neural network (**Wire-Fitted Q-Learning**)
- Cannot just pick the action with the greatest expected utility: must “explore” to be trainable and re-trainable
  - Use a **Boltzmann Probability Distribution** selection policy

## Results

Results were gathered both from simulation and hardware for a variety of tasks. Gathered data for each task included the number of learning iterations, or how many pieces of reward it took for Fido to master the task, action selection time, or the latency in probabalistic action selection, and training time, or the latency in updating Fido’s model.

Table 1: Fido Results in Simulation

Task	Learning Iterations	Action Selection (ms)	Training Time (ms)
Flash	6	0.	6
Float to Point	14	1	6
Drive to Point	17	1	11
Line Following	18	0.	2
Noisy Line Following	21	0.	105

Table 2: Fido Results on Thing One

Task	Learning Iterations	Action Selection (ms)	Training Time (ms)
Stay Still	3	1	43.5
Drive to Point	18	4	65

Table 3: Fido Results on Thing Two

Task	Learning Iterations	Action Selection (ms)	Training Time (ms)
Drive Straight	13	2	30
Line Following	15	21	95
Fetch	8	1	70
Limping Line Following	6	20	37

## Implementation

Fido was programmed in C++, with no external dependencies. However, the simulator does use the SFML graphics library. The hardware implementation uses the Intel Edison embedded platform, a 3D printed chassis and a differential drive system.

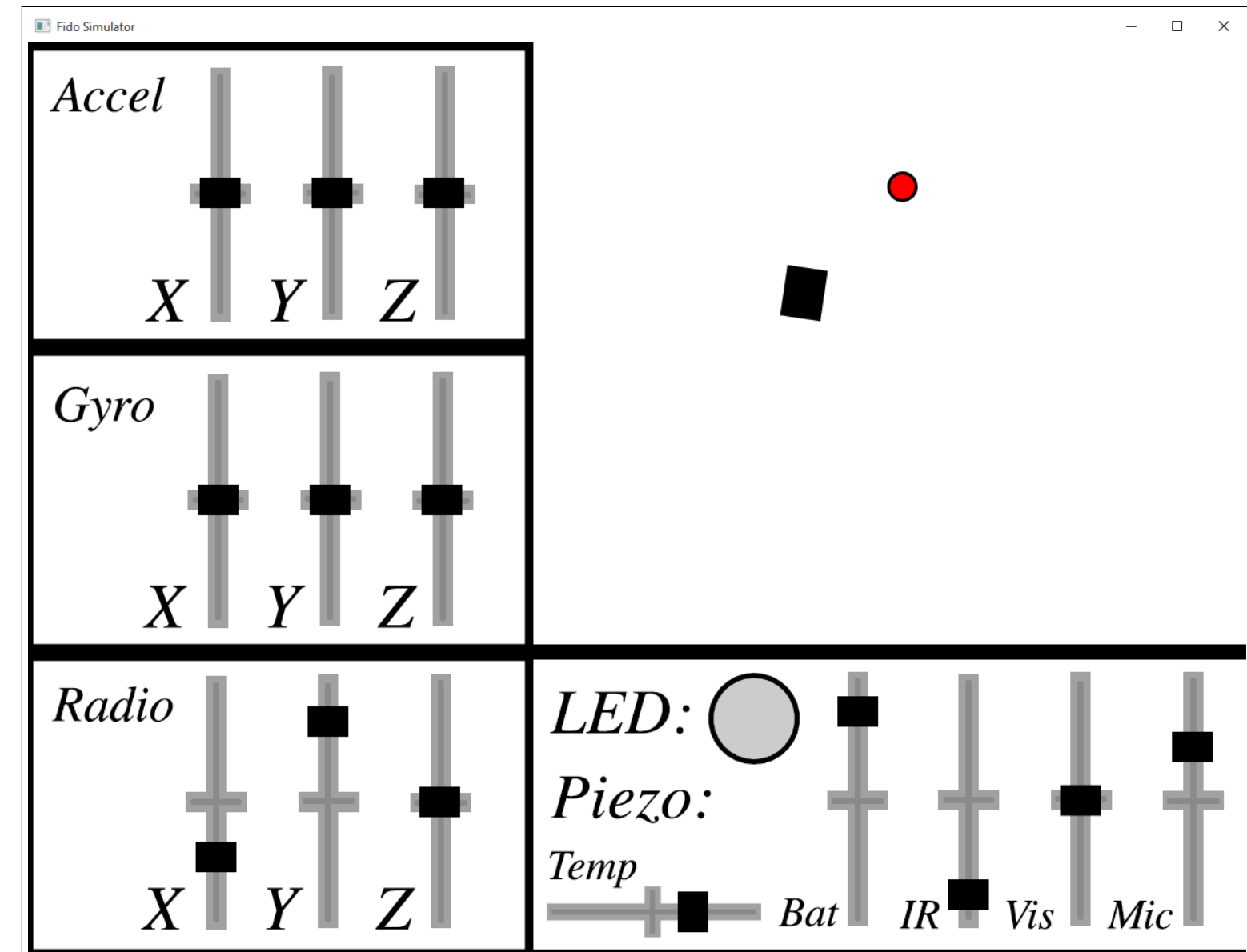


Figure 4: Fido Simulator Graphical User Interface

## Applications

We would like to experiment with dynamic optimization of hyperparameters, changing factors such as neural network architecture and Boltzman temperature constant to best fit the task at hand. We also plan to package Fido as a machine learning library for embedded electronics and robotics, and build a microcontroller-based hardware implementation to further explore for resource-limited environments.

## Selected References

- [1] S. Dini and M. Serrano, “Combining q-learning with artificial neural networks in an adaptive light seeking robot,” 2012.
- [2] C. Gaskett, D. Wettergreen, and A. Zelinsky, “Q-learning in continuous state and action spaces,” pp. 417–428, 1999.
- [3] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” *arXiv preprint arXiv:1312.5602*, 2013.
- [4] M. D. Zeiler, “Adadelta: an adaptive learning rate method,” *arXiv preprint arXiv:1212.5701*, 2012.
- [5] R. Reed, “Pruning algorithms-a survey,” *Neural Networks, IEEE Transactions on*, vol. 4, no. 5, pp. 740–747, 1993.