**DeepRob**

Lecture 14
Imitation Learning I
University of Minnesota

HOW TO TRAIN YOUR DRAGON THE HIDDEN WORLD Kit Harington Auditions with Toothless

# Project 3 — Releases today

- Instructions available on the website
  - Here: https://rpm-lab.github.io/CSCI5980-F24-DeepRob/projects/project3/

  - Uses PROPS Detection dataset

- Implement CNN for classification and Faster R-CNN for detection

- Autograder will be available soon!
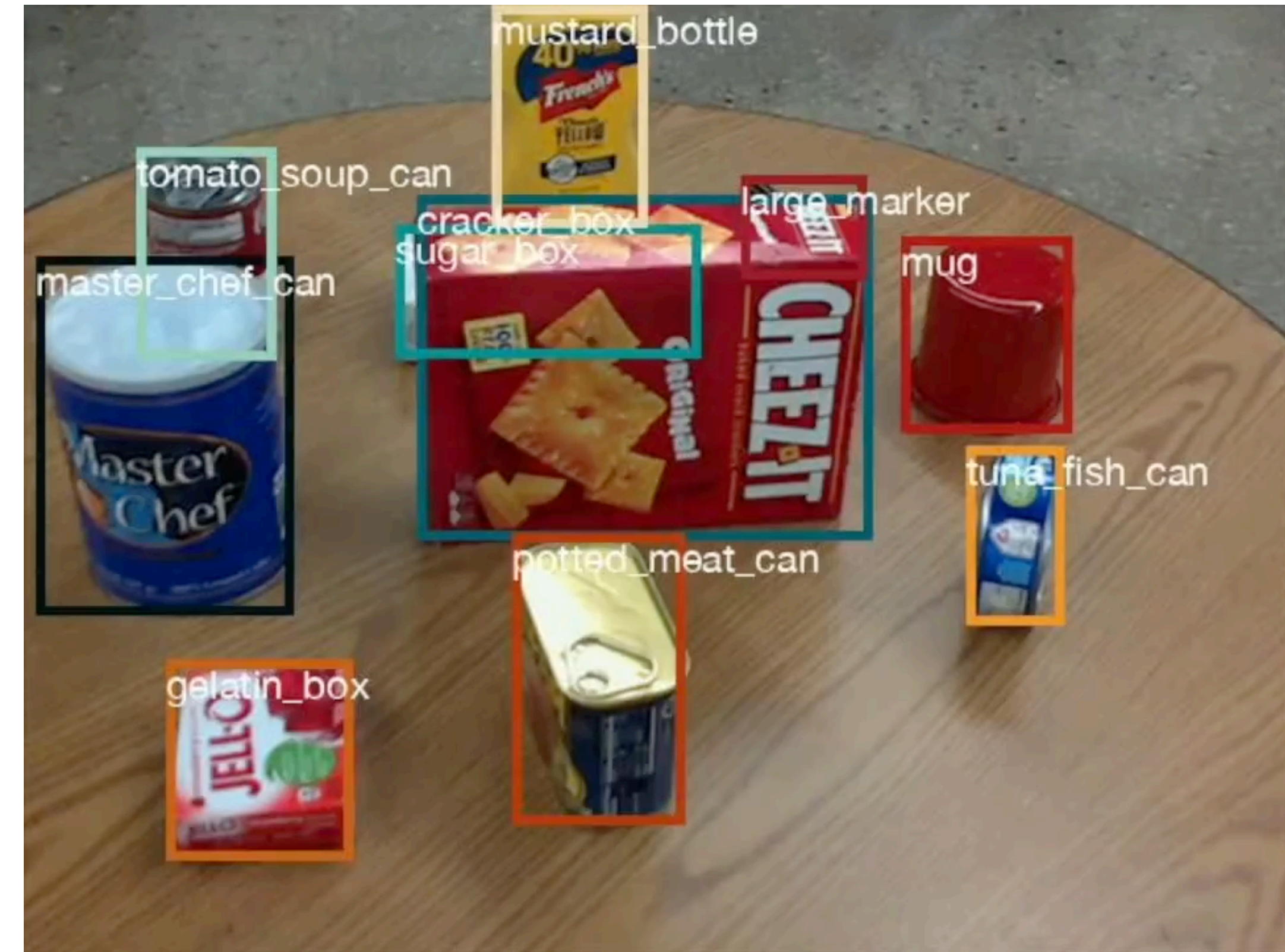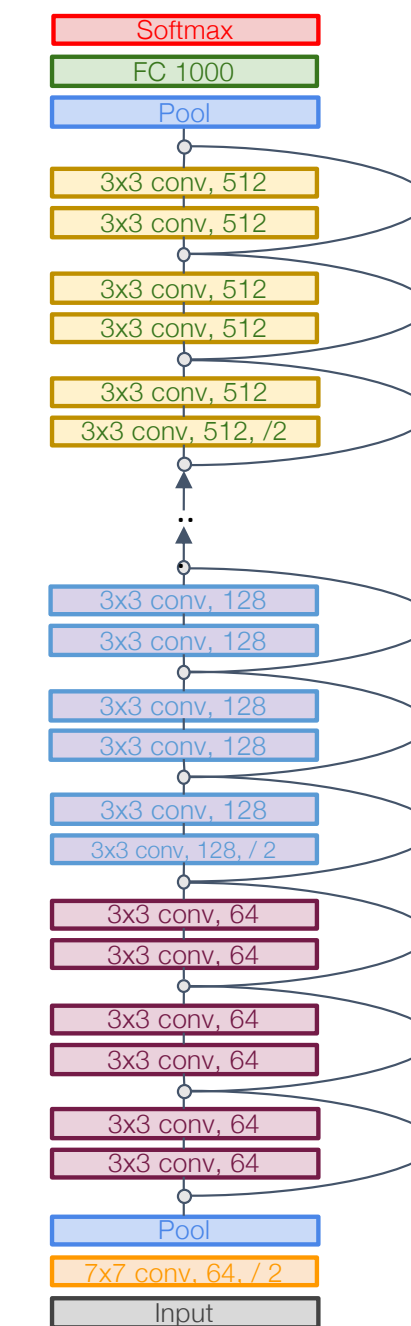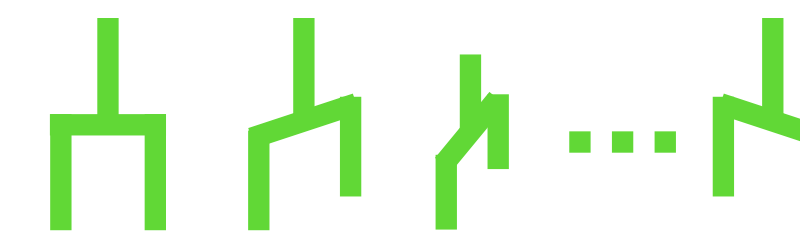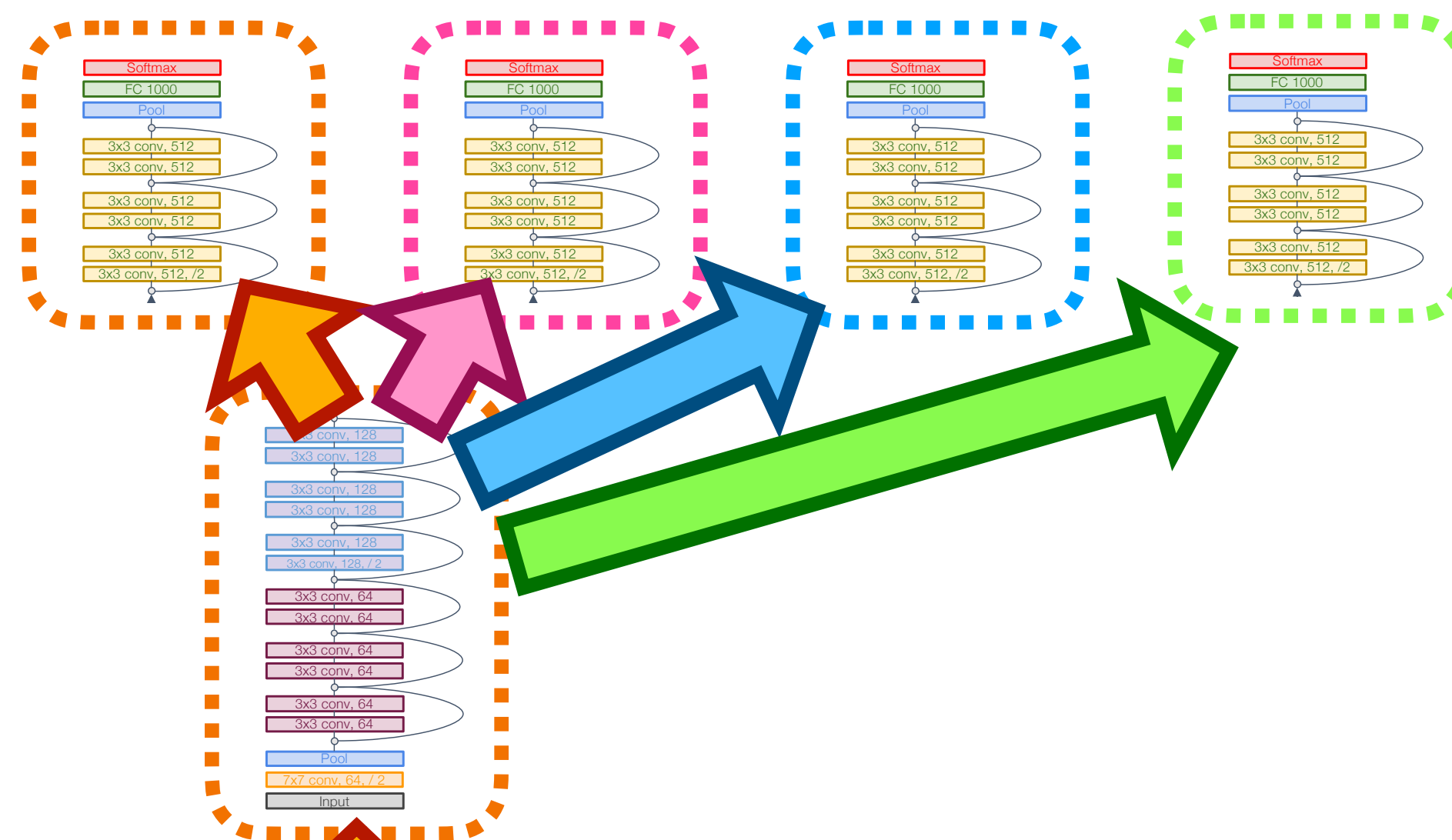
- **Due Monday, October 28th 11:59 PM CT**

**Image Classification**

Cheez It

**Object Detection**

**Segmentation & Pose Estimation**

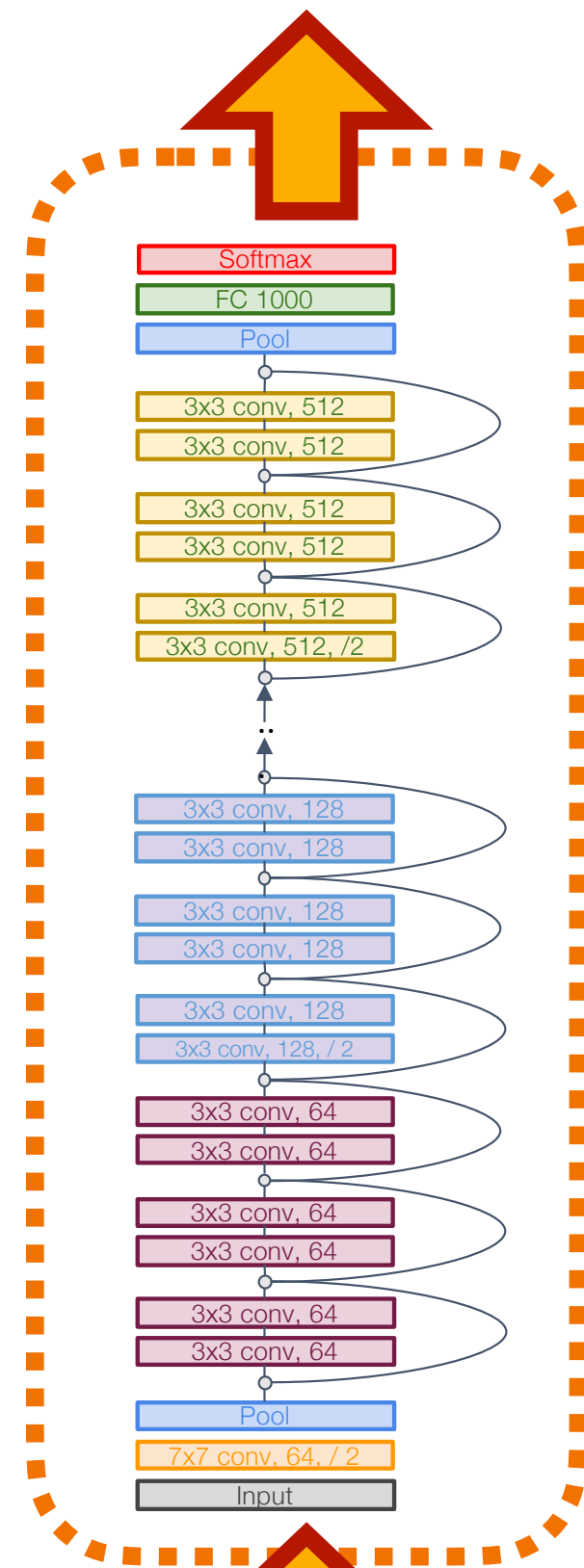**Grasp Detection**

**Grasp Detection**

??? 

???

- After finding grasp poses, how to execute actions?
  - Remember! Inverse Kinematics (5551)
  - Remember! Motion Planning (5551)

**DR**

Action Prediction

Neural Network

# End-to-End Training of Deep Visuomotor Policies

**Sergey Levine**[†]            SVLEVINE@EECS.BERKELEY.EDU
**Chelsea Finn**[†]            CBFINN@EECS.BERKELEY.EDU
**Trevor Darrell**            TREVOR@EECS.BERKELEY.EDU
**Pieter Abbeel**            PABBEEL@EECS.BERKELEY.EDU
*Division of Computer Science*
*University of California*
*Berkeley, CA 94720-1776, USA*
[†]These authors contributed equally.

Action Prediction

# End-to-End Training of Deep Visuomotor Policies

**Sergey Levine**[†]                                    SVLEVINE@EECS.BERKELEY.EDU
**Chelsea Finn**[†]                                      CBFINN@EECS.BERKELEY.EDU
**Trevor Darrell**                                       TREVOR@EECS.BERKELEY.EDU
**Pieter Abbeel**                                        PABBEEL@EECS.BERKELEY.EDU
*Division of Computer Science*
*University of California*
*Berkeley, CA 94720-1776, USA*
[†]These authors contributed equally.
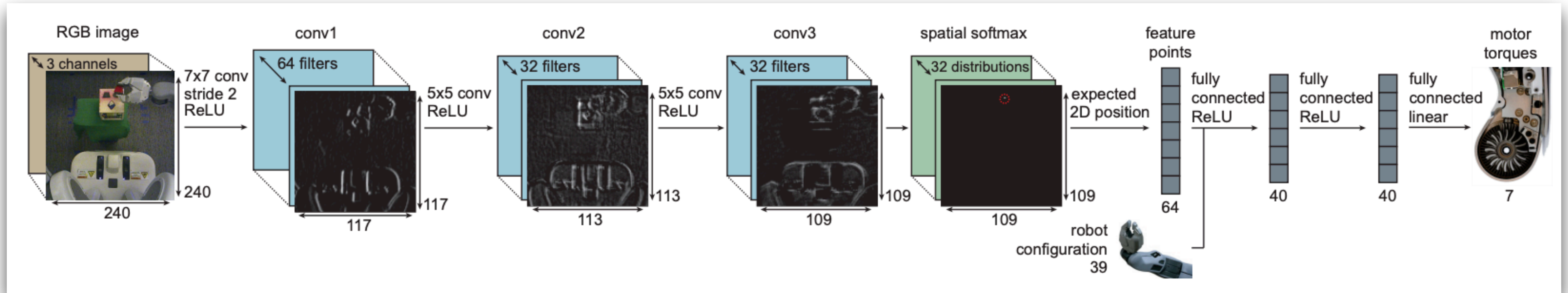
Learning policies that map raw image observations directly to torques at the robot's motors

# What does this entail?



- Input: $o_t$
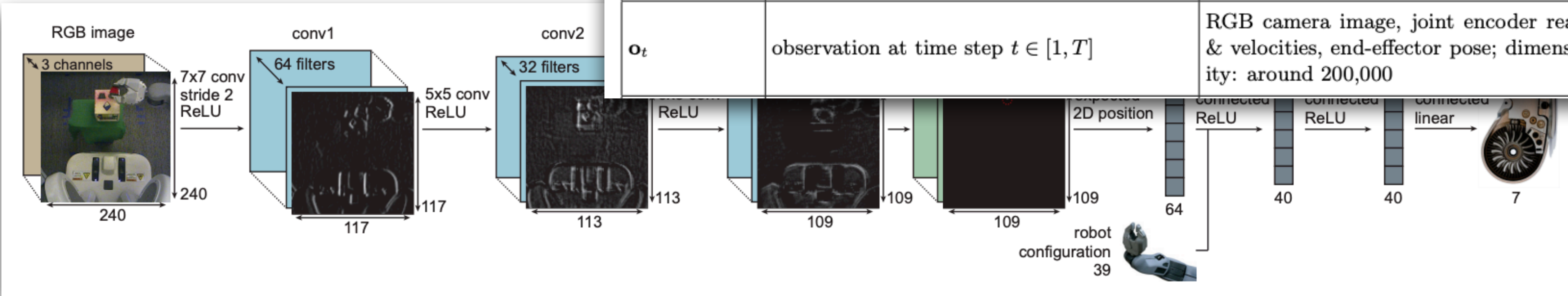- Output: $u_t$
- Policy: $\pi_\theta(u_t \mid o_t)$

Levine, Sergey, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. "End-to-end training of deep visuomotor policies." *Journal of Machine Learning Research* 17, no. 39 (2016): 1-40.

# What c...

| symbol | definition | example/details |
|---|---|---|
| $x_t$ | Markovian system state at time step $t \in [1, T]$ | joint angles, end-effector pose, object positions, and their velocities; dimensionality: 14 to 32 |
| $u_t$ | control or action at time step $t \in [1, T]$ | joint motor torque commands; dimensionality: 7 (for the PR2 robot) |
| $o_t$ | observation at time step $t \in [1, T]$ | RGB camera image, joint encoder readings & velocities, end-effector pose; dimensionality: around 200,000 |



- Input: $o_t$
- Output: $u_t$
- Policy: $\pi_\theta(u_t \mid o_t)$

## State $x_t$
## Vs.
## Observation $o_t$

Levine, Sergey, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. "End-to-end training of deep visuomotor policies." *Journal of Machine Learning Research* 17, no. 39 (2016): 1-40.

RGB image — 3 channels — 240 × 240 — 7x7 conv stride 2 ReLU — conv1 — 64 filters — 117 × 117 — 5x5 conv ReLU — conv2 — 32 filters — 113



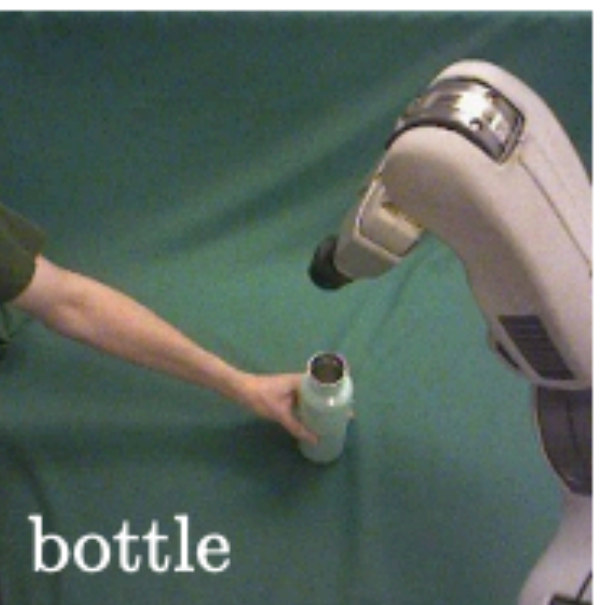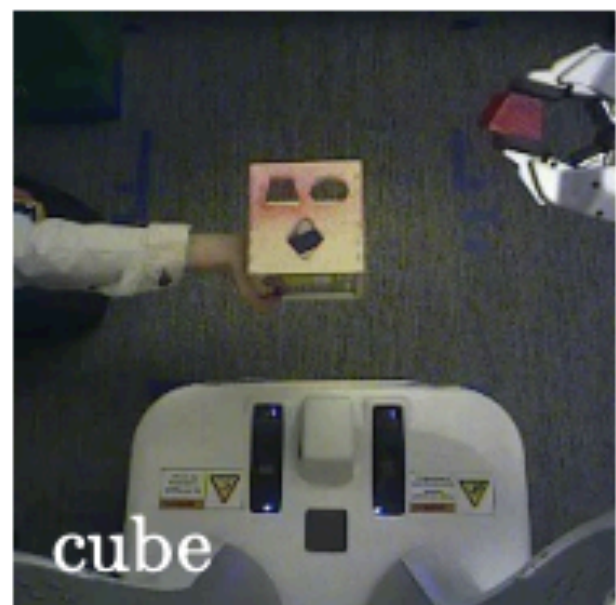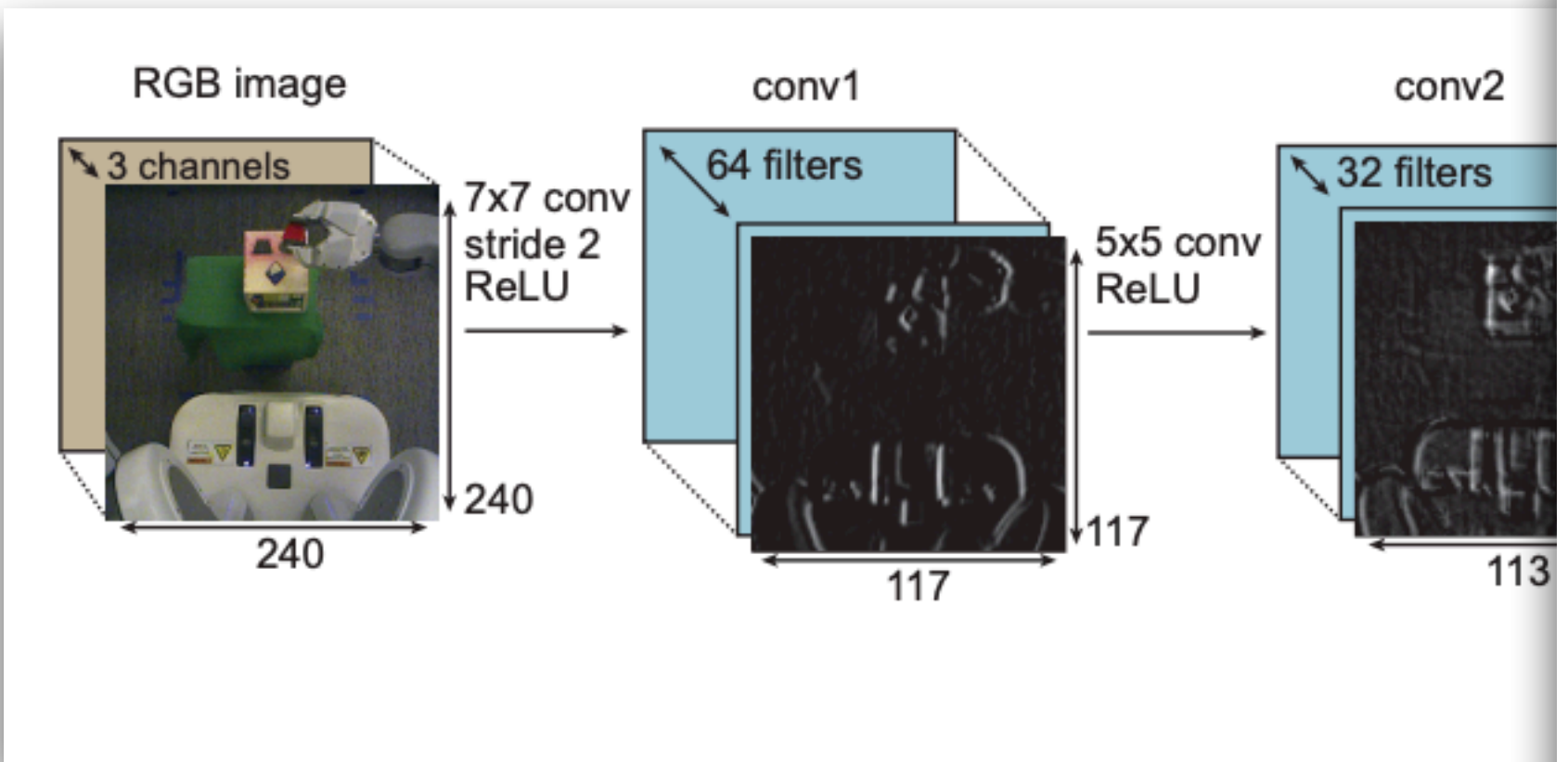hanger   cube   hammer   bottle

Figure 1: Our method learns visuomotor policies that directly use camera image observations (left) to set motor torques on a PR2 robot (right).

| symbol | definition | example/details |
|---|---|---|
| $\mathbf{x}_t$ | Markovian system state at time step $t \in [1, T]$ | joint angles, end-effector pose, object positions, and their velocities; dimensionality: 14 to 32 |
| $\mathbf{u}_t$ | control or action at time step $t \in [1, T]$ | joint motor torque commands; dimensionality: 7 (for the PR2 robot) |
| $\mathbf{o}_t$ | observation at time step $t \in [1, T]$ | RGB camera image, joint encoder readings & velocities, end-effector pose; dimensionality: around 200,000 |
| $\tau$ | trajectory: $\tau = \{\mathbf{x}_1, \mathbf{u}_1, \mathbf{x}_2, \mathbf{u}_2, \ldots, \mathbf{x}_T, \mathbf{u}_T\}$ | notational shorthand for a sequence of states and actions |
| $\ell(\mathbf{x}_t, \mathbf{x}_t)$ | cost function that defines the goal of the task | distance between an object in the gripper and the target |
| $p(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{u}_t)$ | unknown system dynamics | physics that govern the robot and any objects it interacts with |
| | bution | stochastic process that produces camera images from system state |
| | olicy parameter- | convolutional neural network, such as the one in Figure 2 |
| | | notational shorthand for observation-based policy conditioned on state |
| | linear-Gaussian | time-varying linear-Gaussian controller has form $\mathcal{N}(\mathbf{K}_{ti}\mathbf{x}_t + \mathbf{k}_{ti}, \mathbf{C}_{ti})$ |
| | for $\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)$: $\mathbf{x}_t, \mathbf{u}_t)$ | notational shorthand for trajectory distribution induced by a policy |

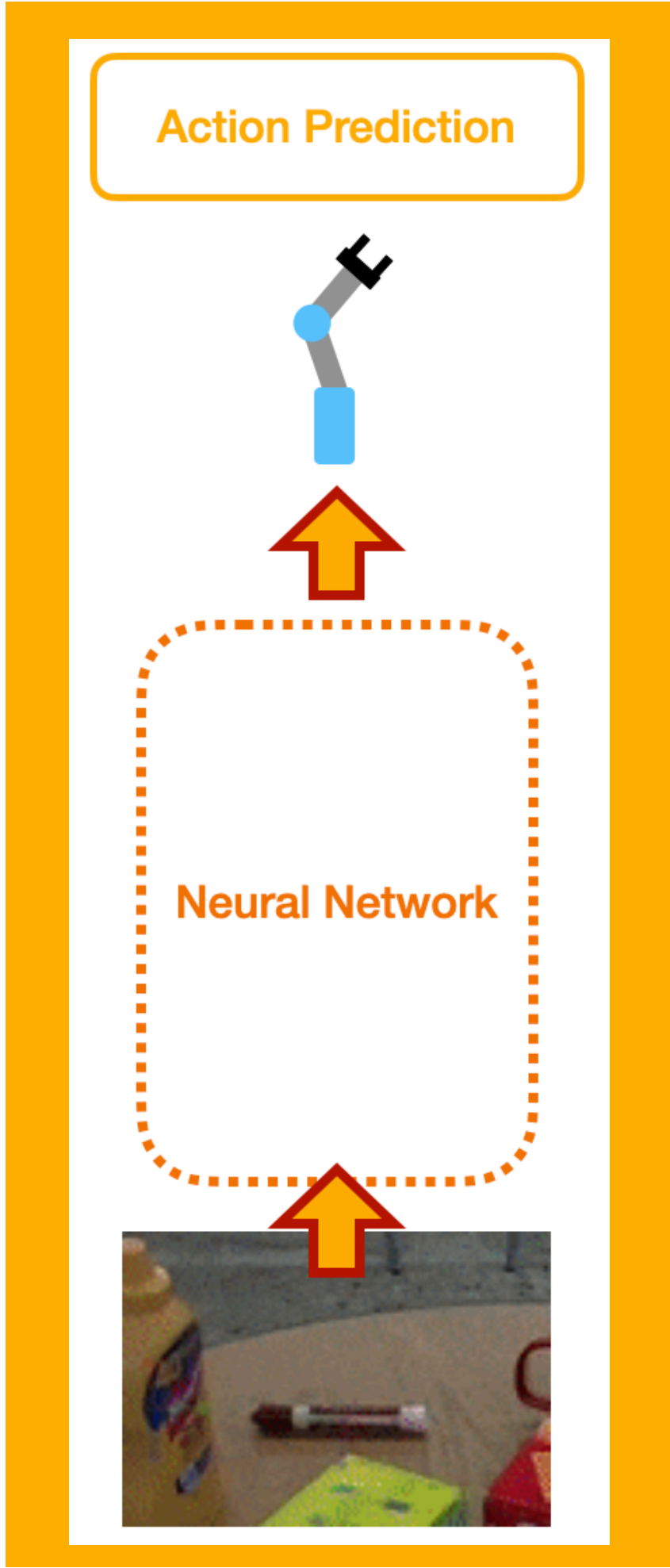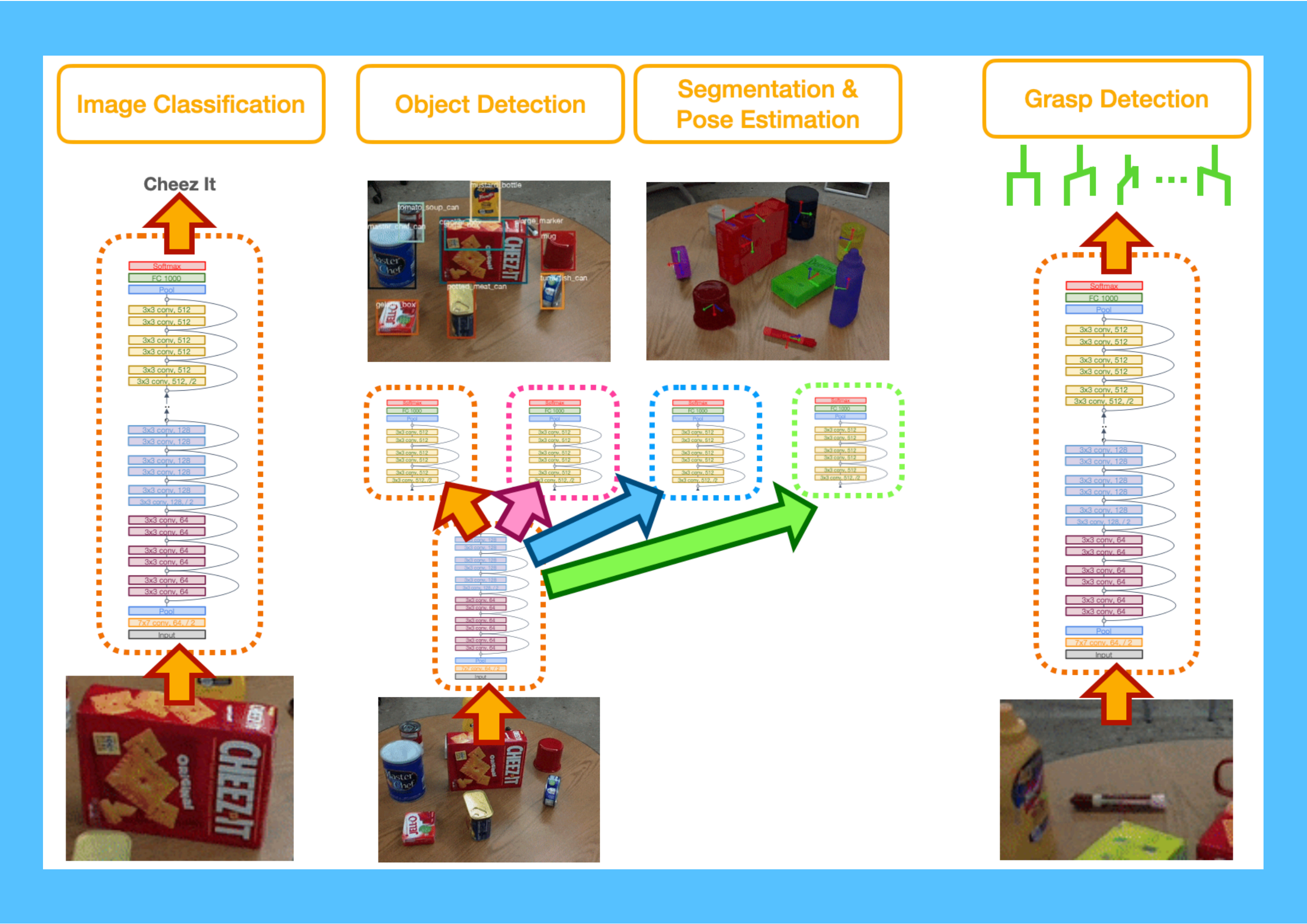Table 1: Summary of the notation frequently used in this article.

Levine, Sergey, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. "End-to-end training of deep visuomotor policies." *Journal of Machine Learning Research* 17, no. 39 (2016): 1-40.

Levine, Sergey, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. "End-to-end training of deep visuomotor policies." *Journal of Machine Learning Research* 17, no. 39 (2016): 1-40.
https://www.youtube.com/watch?v=Q4bMcUk6pcw

# Challenges in going from **Prediction** to **Control**

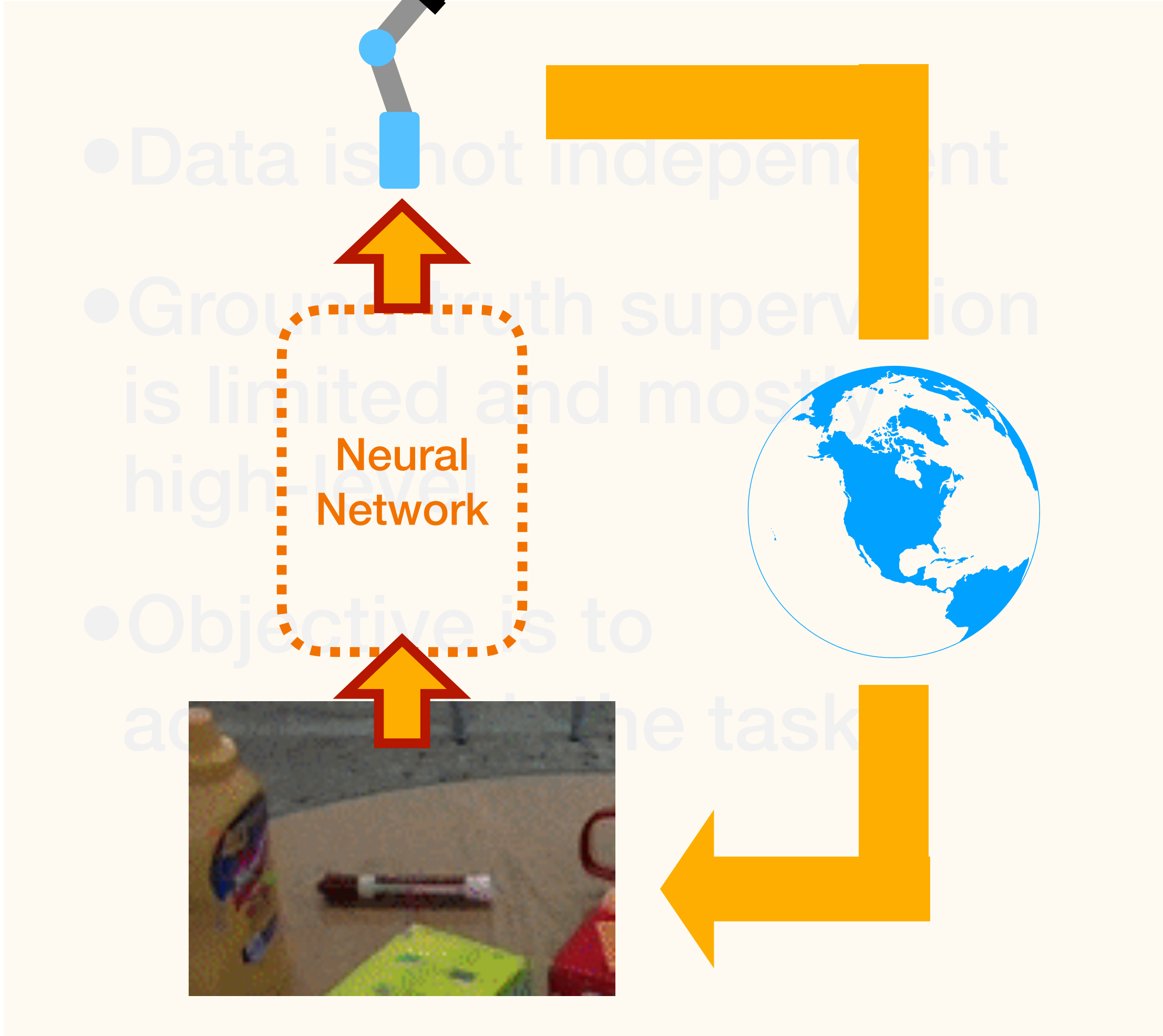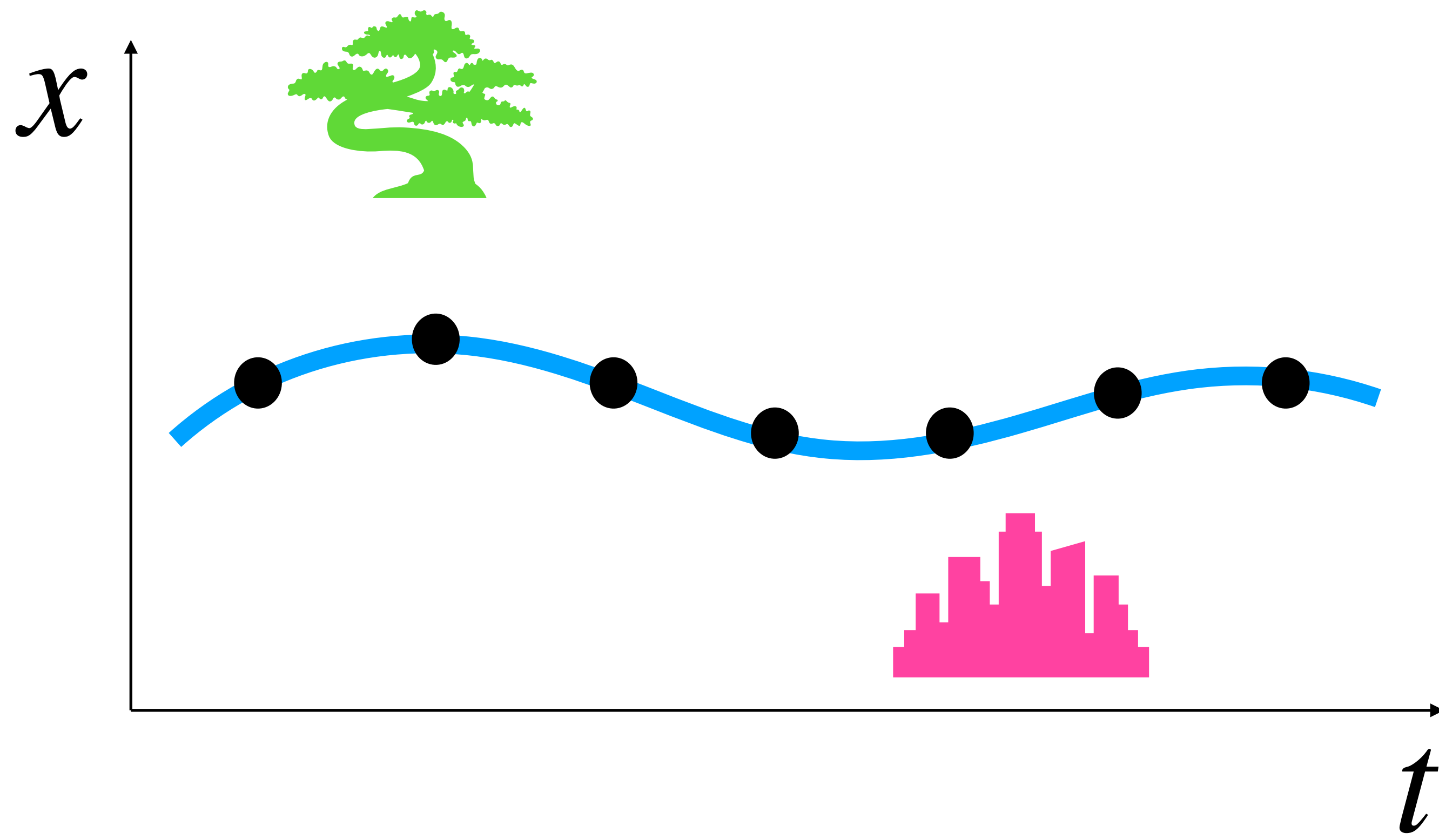# Challenges in going from **Prediction** to **Control**

- Data is i.i.d distributed

- Ground truth supervision for the prediction is available

- Objective is to predict the right label or regress a value close to the ground truth

- Data is not independent

- Ground truth supervision is limited and mostly high-level

- Objective is to accomplish the task

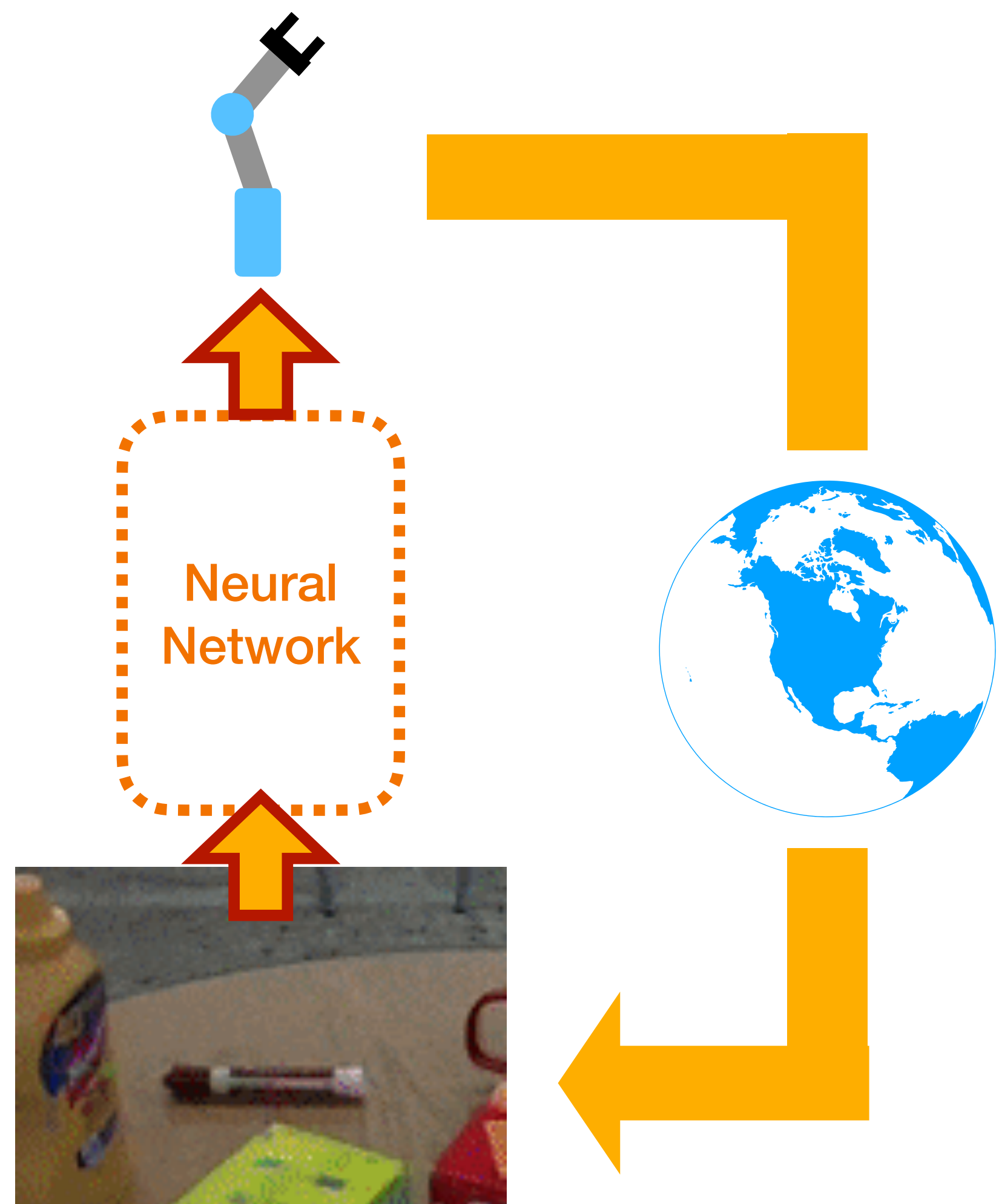# Challenges in going from **Prediction** to **Control**

- Data is i.i.d distributed

- Ground truth supervision for the prediction is available

- Objective is to predict the right label or regress a value close to the ground truth

- Data is not independent

- Ground truth supervision is limited and mostly high-level

- Objective is to achieve the task

Neural Network

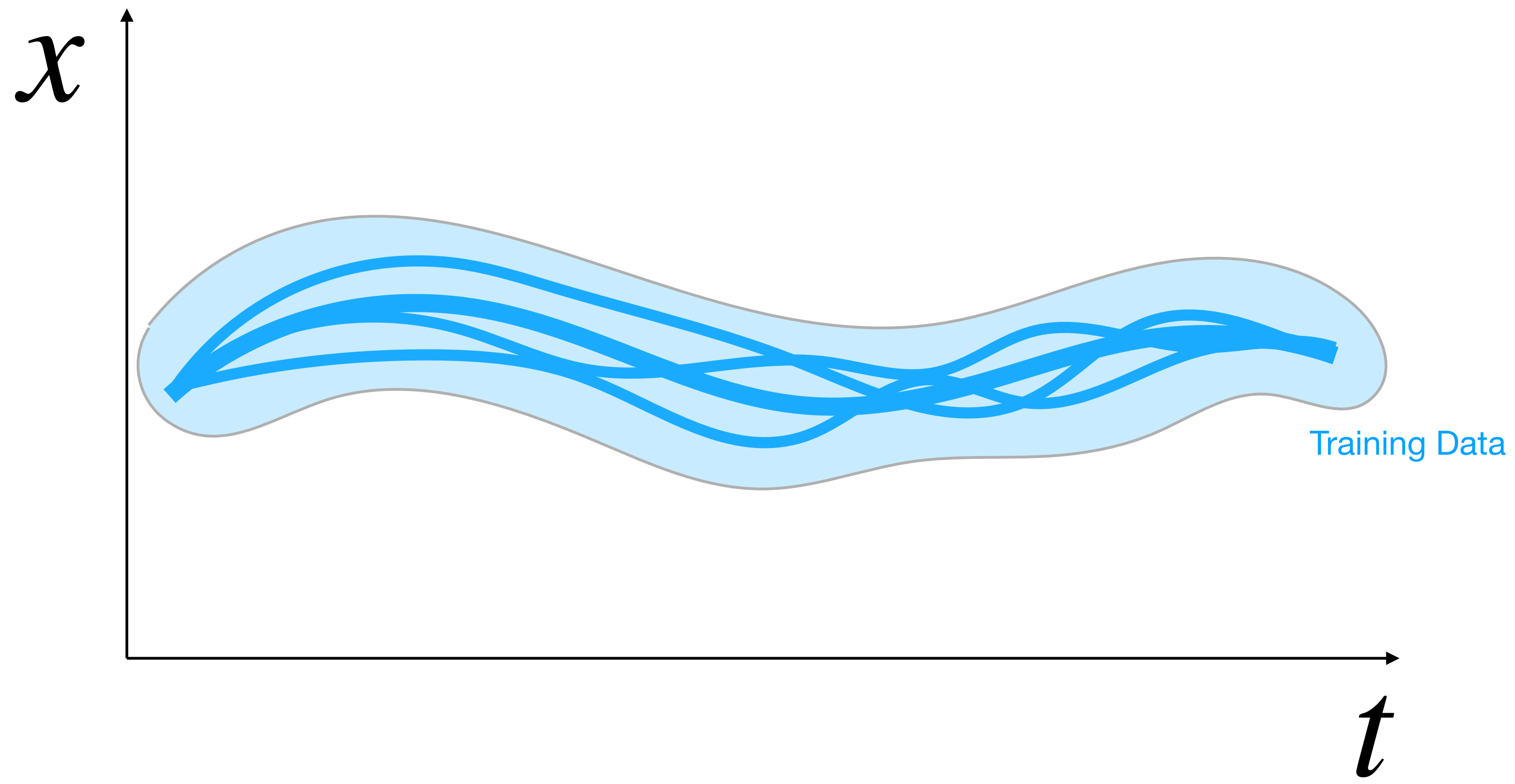There is feedback and associated issues!

# There is feedback and associated issues!



Data is dependent

# There is feedback and associated issues!
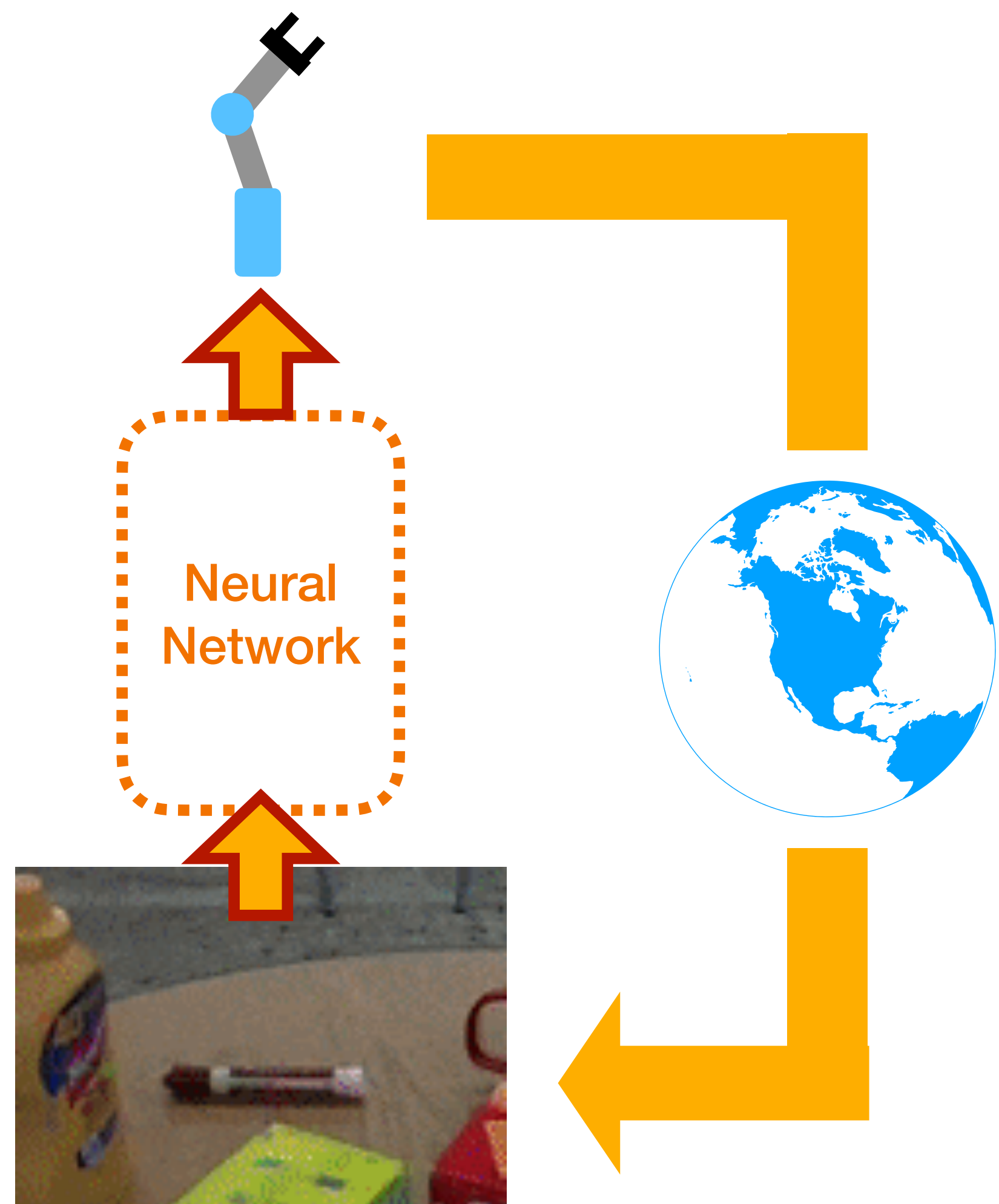


Training Data

# There is feedback and associated issues!



Learned Data

Training Data

$x$

$t$

Neural Network

Covariance Shift due to Feedback

# This is a commonly seen issue

## Exploring the Limitations of Behavior Cloning for Autonomous Driving

Felipe Codevilla *
Computer Vision Center (CVC)
Campus UAB, Barcelona, Spain
fcodevilla@cvc.uab.es

Eder Santana
Toyota Research Institute (TRI)
Los Altos, CA, USA.
edercsjr@gmail.com

Antonio M. López
Computer Vision Center (CVC)
Campus UAB, Barcelona, Spain
antonio@cvc.uab.es

Adrien Gaidon
Toyota Research Institute (TRI)
Los Altos, CA, USA.
adrien.gaidon@tri.global

## ChauffeurNet: Learning to Drive by Imitating the Best and Synthesizing the Worst

Mayank Bansal
Waymo Research
Mountain View, CA, USA
mayban@waymo.com

Alex Krizhevsky[†]
akrizhevsky@gmail.com

Abhijit Ogale
Waymo Research
Mountain View, CA, USA
ogale@waymo.com

## Imitating Driver Behavior with Generative Adversarial Networks

Alex Kuefler[1], Jeremy Morton[2], Tim Wheeler[2], and Mykel Kochenderfer[2]

## Causal Confusion in Imitation Learning

Pim de Haan[*,1], Dinesh Jayaraman[†,‡], Sergey Levine[†]
*Qualcomm AI Research, University of Amsterdam,
[†]Berkeley AI Research, [‡] Facebook AI Research

# DAGGER

## A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning

| Stéphane Ross | Geoffrey J. Gordon | J. Andrew Bagnell |
|---|---|---|
| Robotics Institute | Machine Learning Department | Robotics Institute |
| Carnegie Mellon University | Carnegie Mellon University | Carnegie Mellon University |
| Pittsburgh, PA 15213, USA | Pittsburgh, PA 15213, USA | Pittsburgh, PA 15213, USA |
| stephaneross@cmu.edu | ggordon@cs.cmu.edu | dbagnell@ri.cmu.edu |

Initialize $\mathcal{D} \leftarrow \emptyset$.
Initialize $\hat{\pi}_1$ to any policy in $\Pi$.
**for** $i = 1$ **to** $N$ **do**
  Let $\pi_i = \beta_i \pi^* + (1 - \beta_i)\hat{\pi}_i$.
  Sample $T$-step trajectories using $\pi_i$.
  Get dataset $\mathcal{D}_i = \{(s, \pi^*(s))\}$ of visited states by $\pi_i$
  and actions given by expert.
  Aggregate datasets: $\mathcal{D} \leftarrow \mathcal{D} \bigcup \mathcal{D}_i$.
  Train classifier $\hat{\pi}_{i+1}$ on $\mathcal{D}$.
**end for**
**Return** best $\hat{\pi}_i$ on validation.

**Algorithm 3.1:** DAGGER Algorithm.

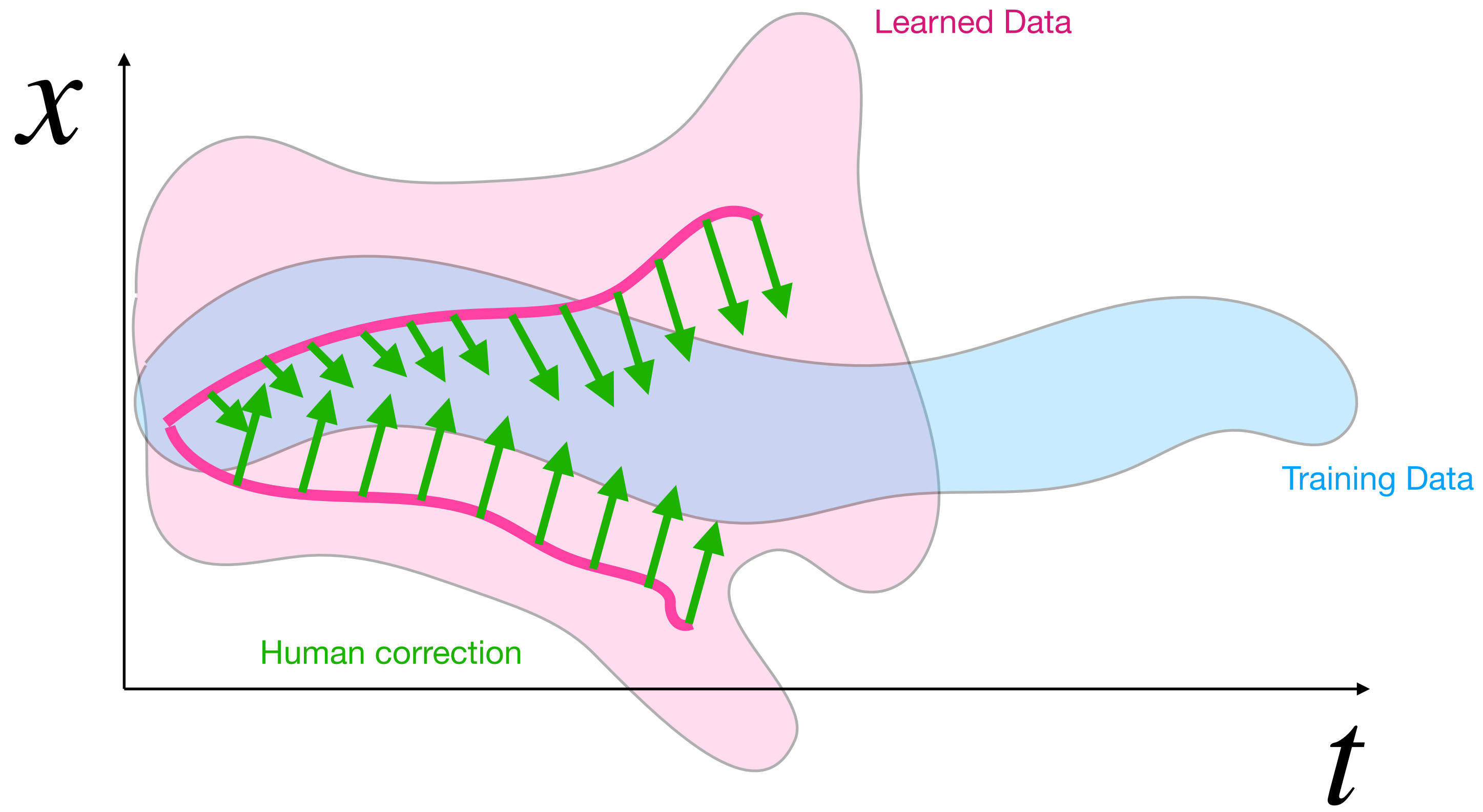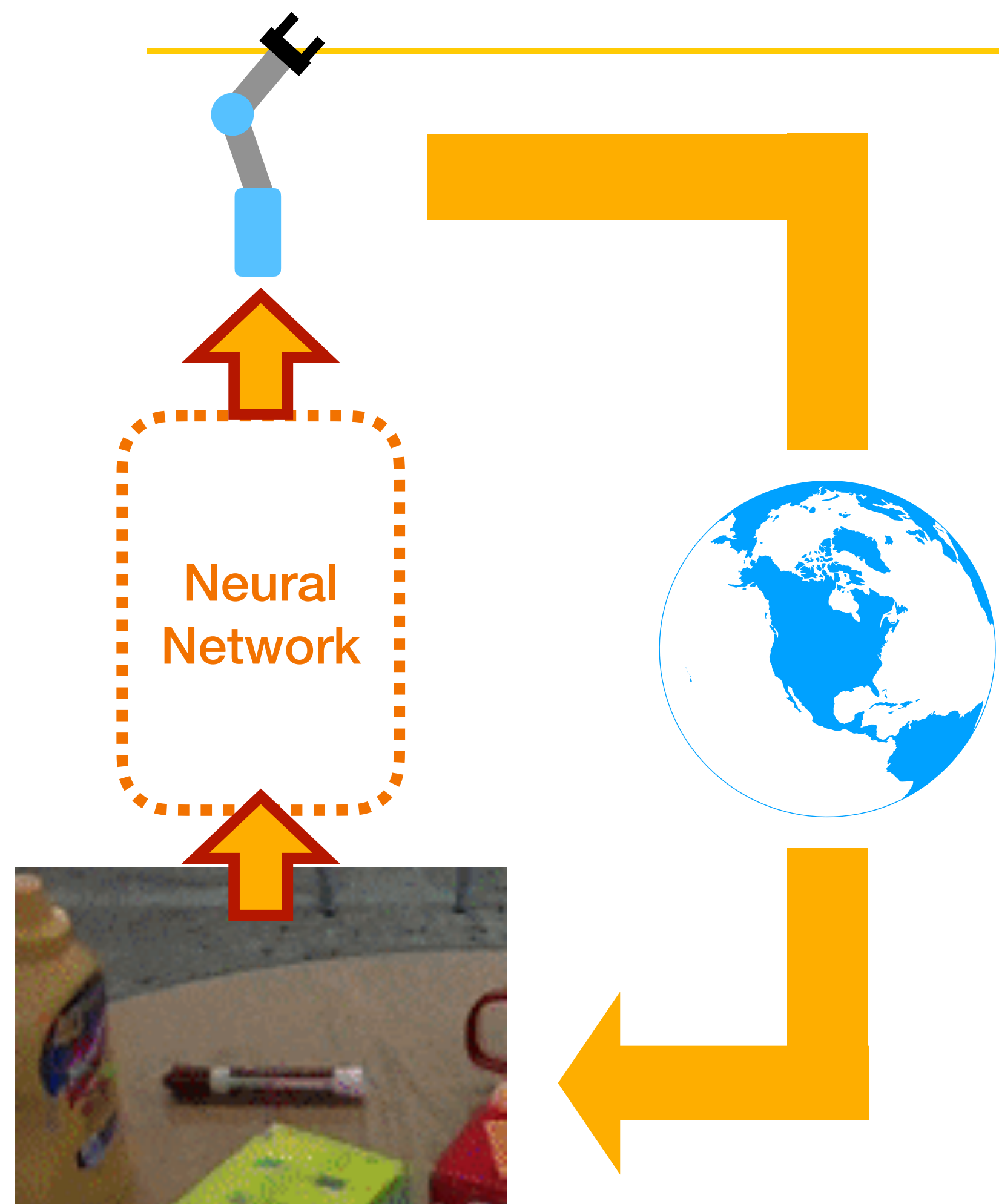Step 1: Collect Human Demonstrations and Train initial policy $\pi$

Step 2: Rollout $\pi(\,.\,)$ to collect new states $x_t$ or observations $o_t$
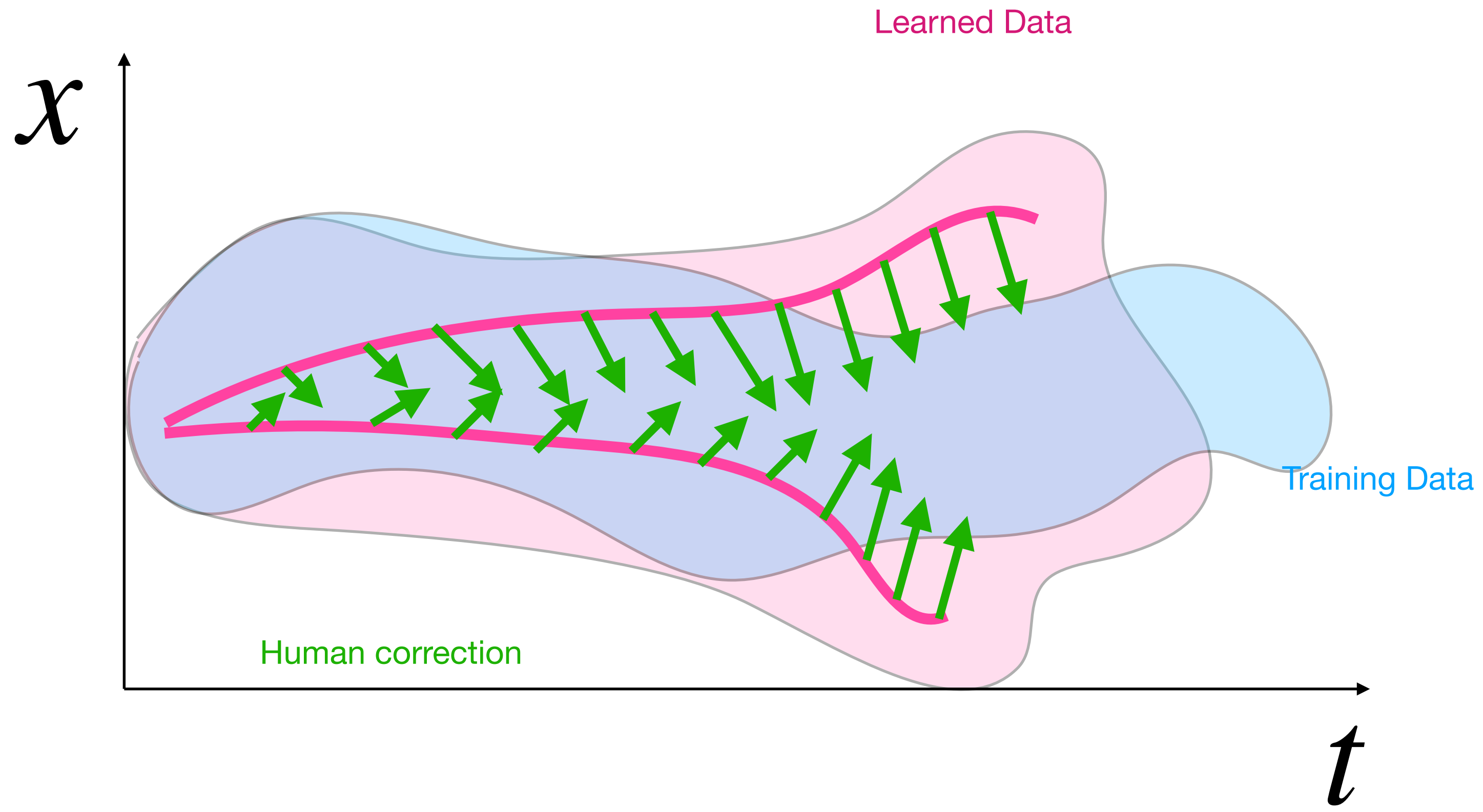
Step 3: Ask human for correct action

Step 4: Aggregate data & train $\pi(\,.\,)$
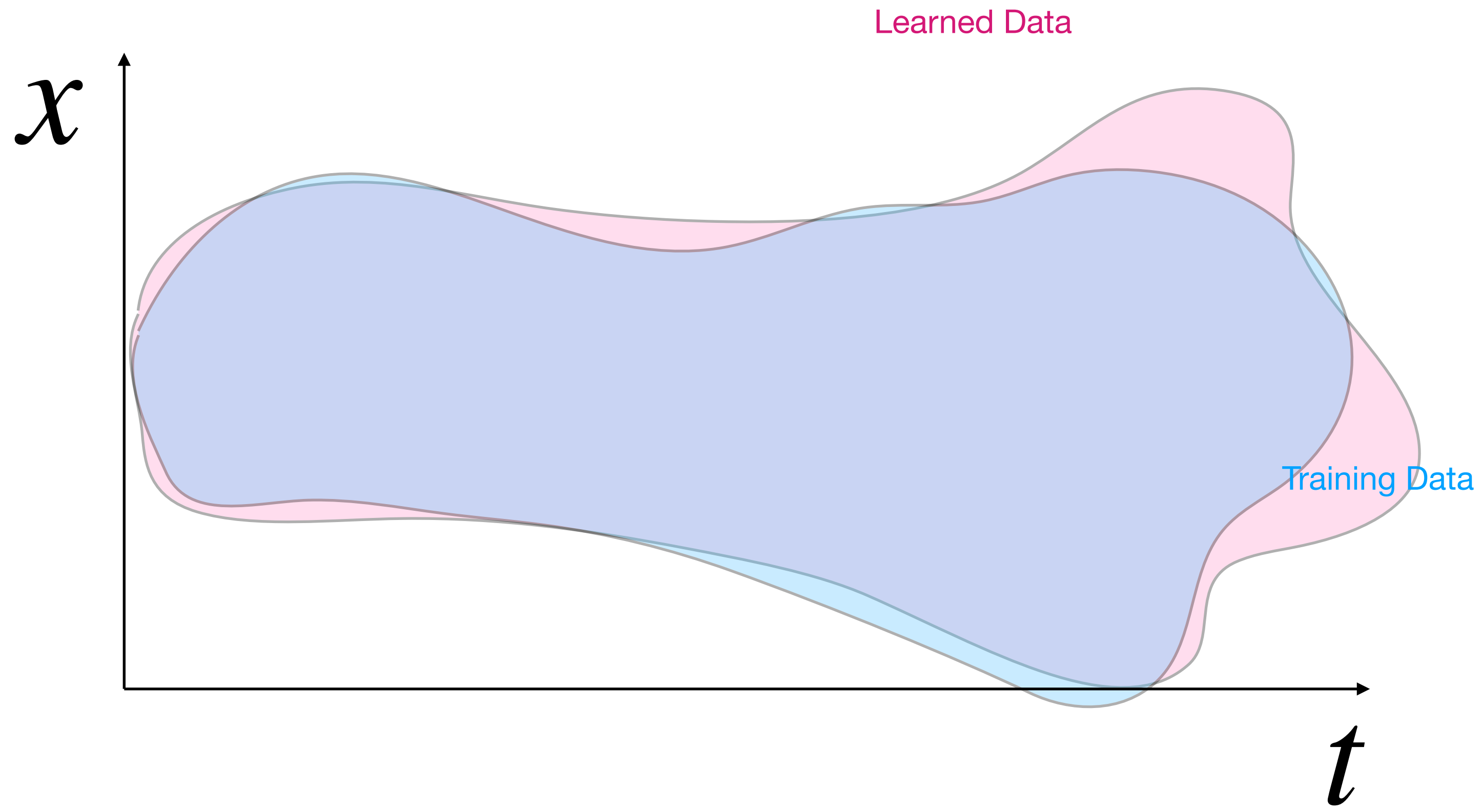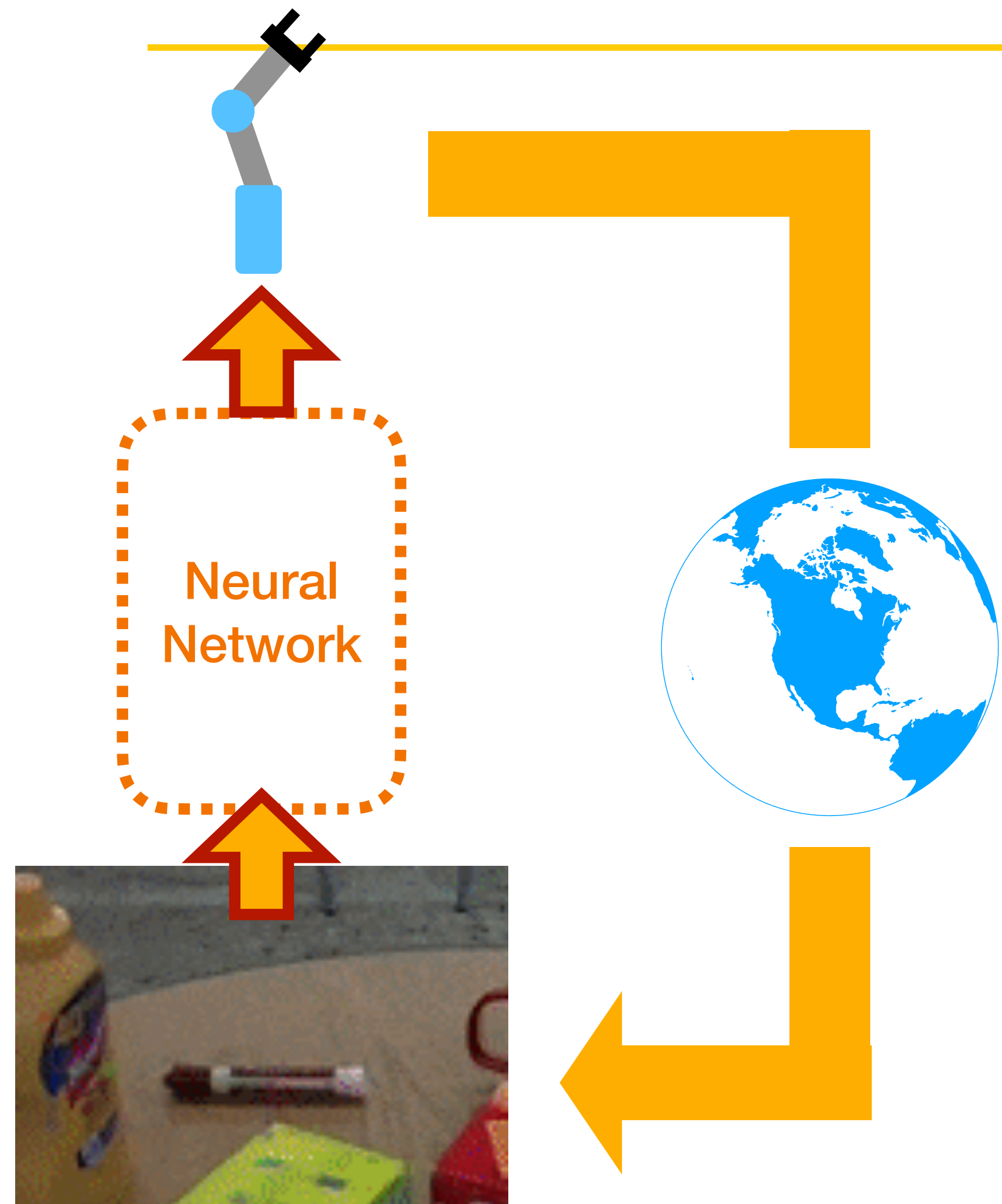
# There is feedback and associated issues!



Learned Data

$x$

Training Data

Human correction

$t$

Neural
Network

Neural
Network

Learned Data

$x$

Training Data

Human correction

$t$

# There is feedback and associated issues!



Learned Data

Training Data

$x$

$t$

# Next Lecture:
# Imitation Learning II

# DeepRob

**Lecture 14**
**Imitation Learning I**
**University of Minnesota**