

# CS 277: Control and Reinforcement Learning

Winter 2021

## Lecture 17: Multi-Task Learning

Roy Fox

Department of Computer Science

Bren School of Information and Computer Sciences

University of California, Irvine



# Today's lecture

---

Transfer learning

Domain randomization + adaptation

Shared learning

# Learning from very little data

---

- As the number of **learnable tasks** grows
  - **sample complexity** per task must drop to be practical
- **Our goal:** learn a new task with
  - **0-shot:** no new training interactions (exploration / demonstration)
  - **1-shot:** single training episode
  - **few-shot:** very few training episodes

# Prior knowledge

---

- To only need little information from data, the rest must be **a-priori**
- **Programmed** prior knowledge:
  - Programmed policy / skills
  - Choice of observation and action **representations**
    - **Feature extraction**
  - **World model** (dynamics / reward)
  - Learner **model class** / neural network **architecture**



# Learning prior knowledge from other tasks

---

- **Transfer learning**: first learn other task(s), then solve new task
  - with ( $>0$ -shot) or without (0-shot) more learning in the new task
- Practical question: what knowledge is **transferred / shared**?
  - **Value function / policy**
  - **Perceptual features**
  - **World model**
  - More later...

# Idea 1: policy transfer

---

- Find **similar task(s)** where data is abundant
  - Many **demonstrations** / exploration **episodes** can be obtained
  - E.g. simulator of the world → real world (**sim2real**)
- **Train** policy with RL / IL in the abundant domain
- **Execute** policy in the scarce domain
  - Or **fine-tune** with further few-shot RL / IL, as needed

# Soft-optimal policies for fine-tuning

---

- **Problem:** policy can “overfit” to pre-training task
  - Policy may become deterministic
    - unfit for exploration
    - optimizer may struggle to again introduce uncertainty
  - Perceptual features may deteriorate to only what's needed for actions
- **Solution:** keep policy soft-optimal
  - Max entropy subject to sufficiently high value

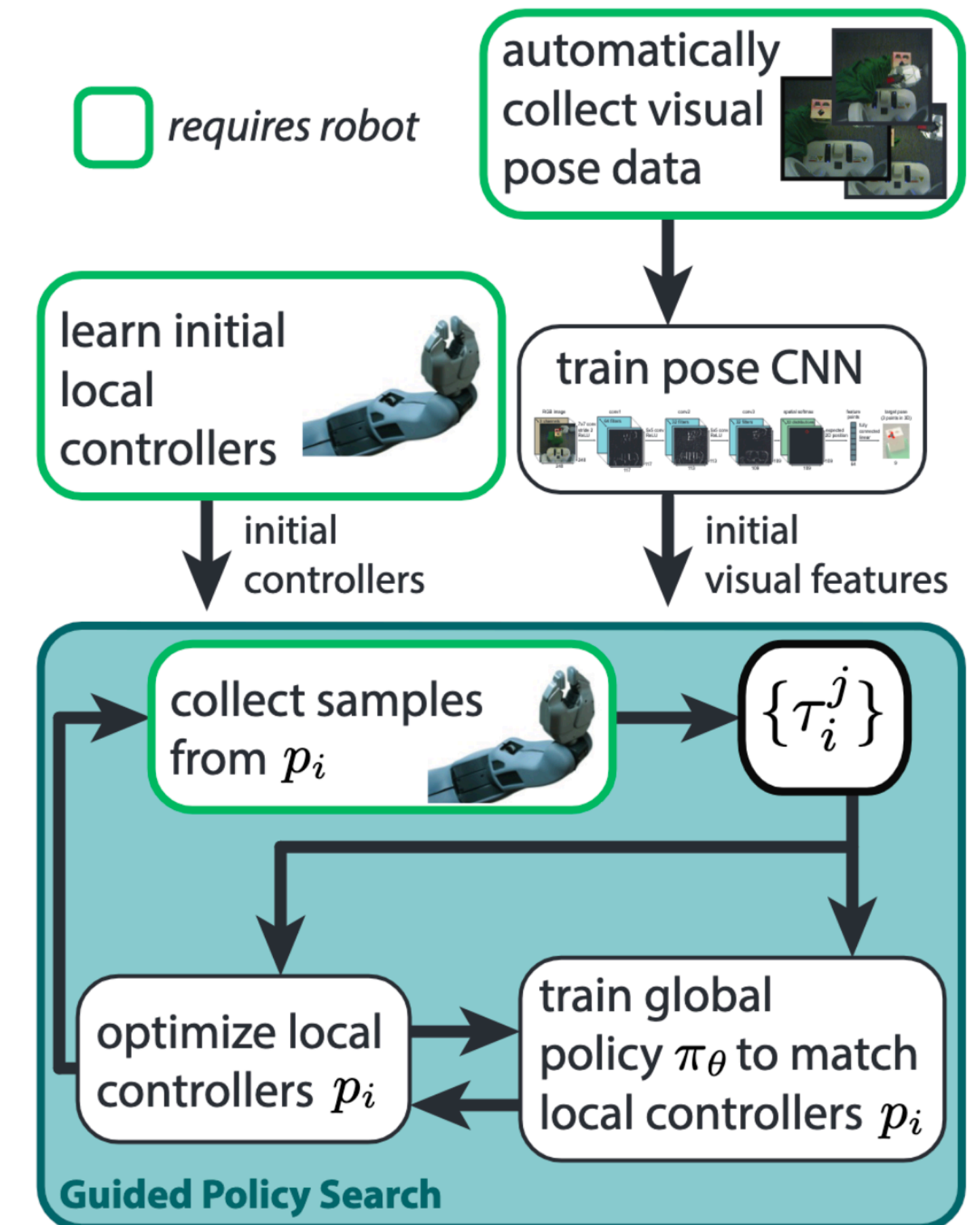
# SQL pre-training helps fine-tuning

---

**Soft Q-learning**  
Fine-tuning a pretrained policy  
in a new environment

# Idea 2: perceptual features transfer

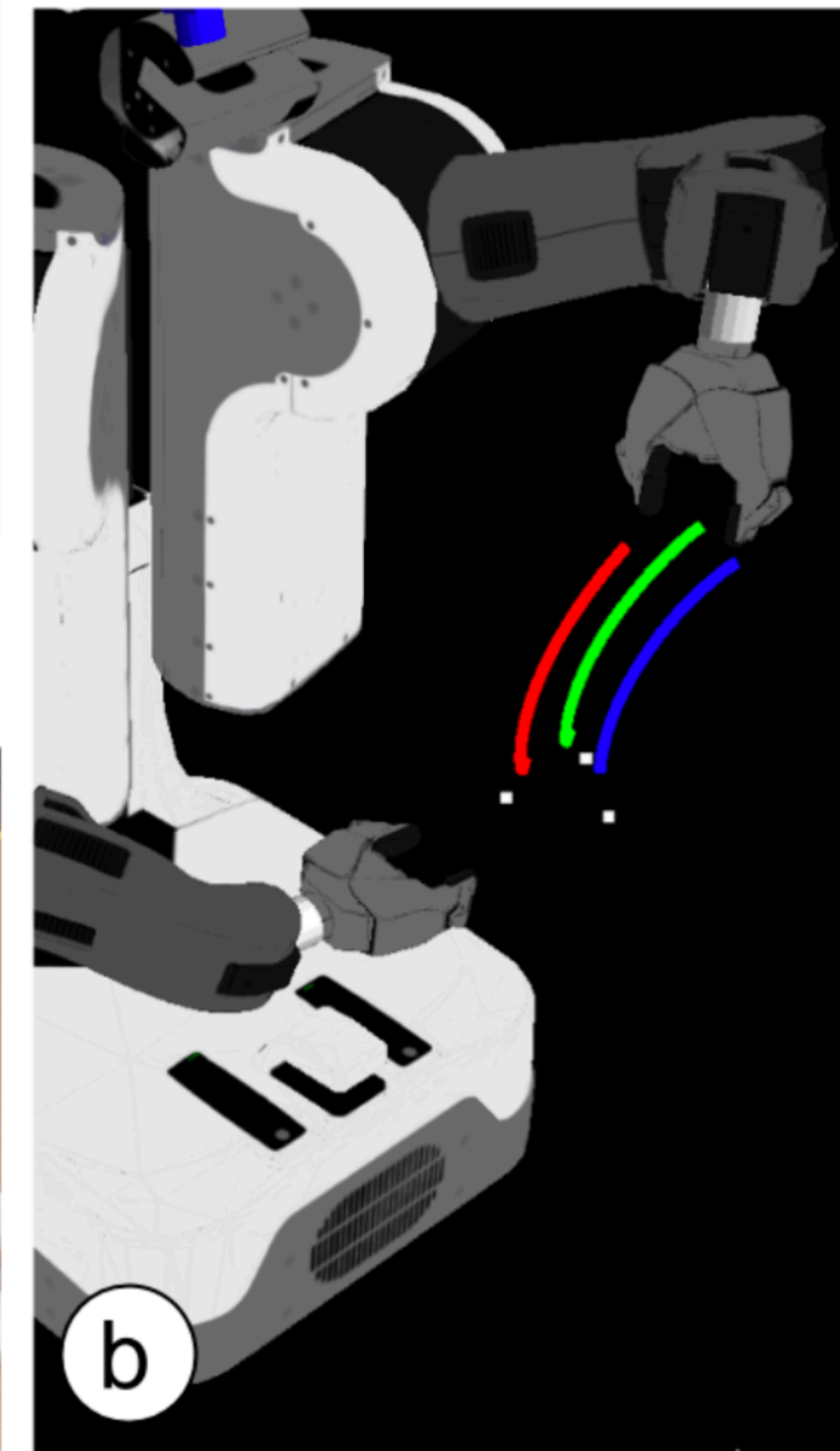
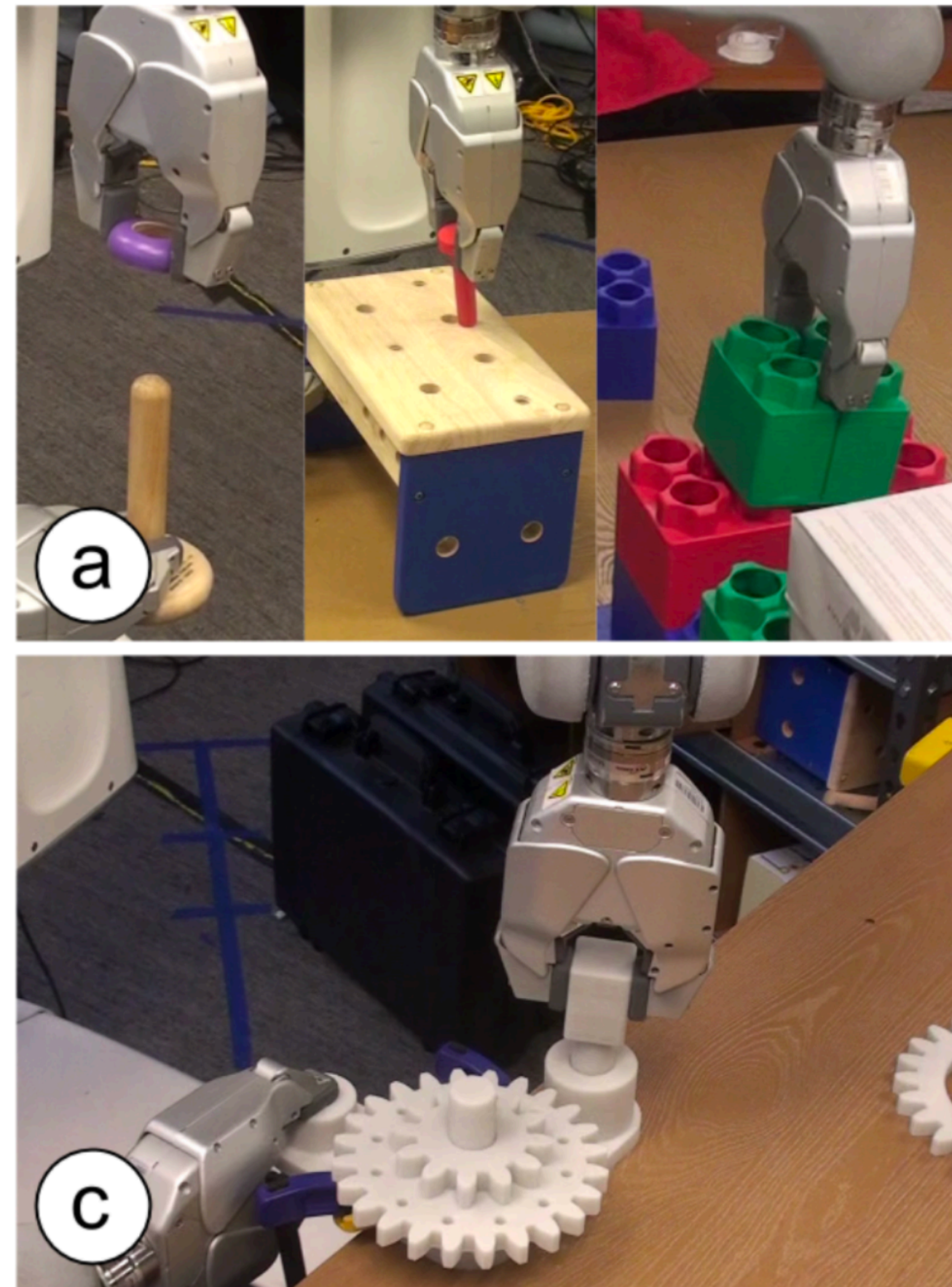
- **Interact** to collect  $\langle$ robot pose, image $\rangle$  data
- Train perceptual **features** that recover image  $\rightarrow$  pose
- **RL** with perceptual features as observations
- May again benefit from **fine-tuning**





# Idea 3: model transfer

- Interact on one / many **related tasks**
- Fit a **world model** to the dynamics
- **Model-based RL** of a new task
- **Problem**: prior model is inaccurate
- **Solution**: take the **pre-trained** model
  - and **fine-tune** it using new task data



# Today's lecture

---

Transfer learning

**Domain randomization + adaptation**

Shared learning

# Domain randomization

---

- Choosing a **source domain** to match the target domain may be hard
- Can we do better with **multiple** source domains?
  - Define **distribution over tasks** that supports the target  $\implies$  **interpolation**
  - **Generalize** even outside the support  $\implies$  **extrapolation**



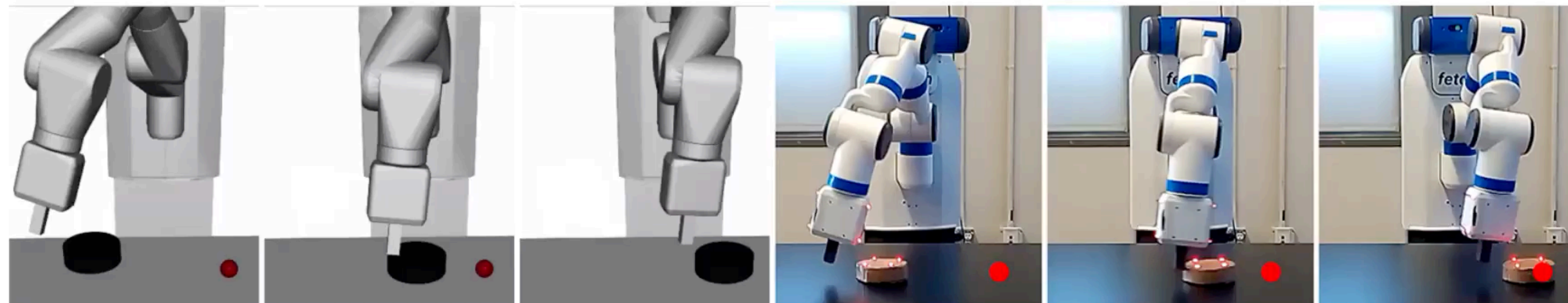
# Sim2real with domain randomization

## Sim-to-Real Transfer of Robotic Control with Dynamics Randomization

Xue Bin Peng<sup>1,2</sup>, Marcin Andrychowicz<sup>2</sup>, Wojciech Zaremba<sup>2</sup>, Pieter Abbeel<sup>1,2</sup>

<sup>1</sup>Electrical Engineering and Computer Sciences, UC Berkeley, USA

<sup>2</sup>OpenAI, USA



# Domain adaptation

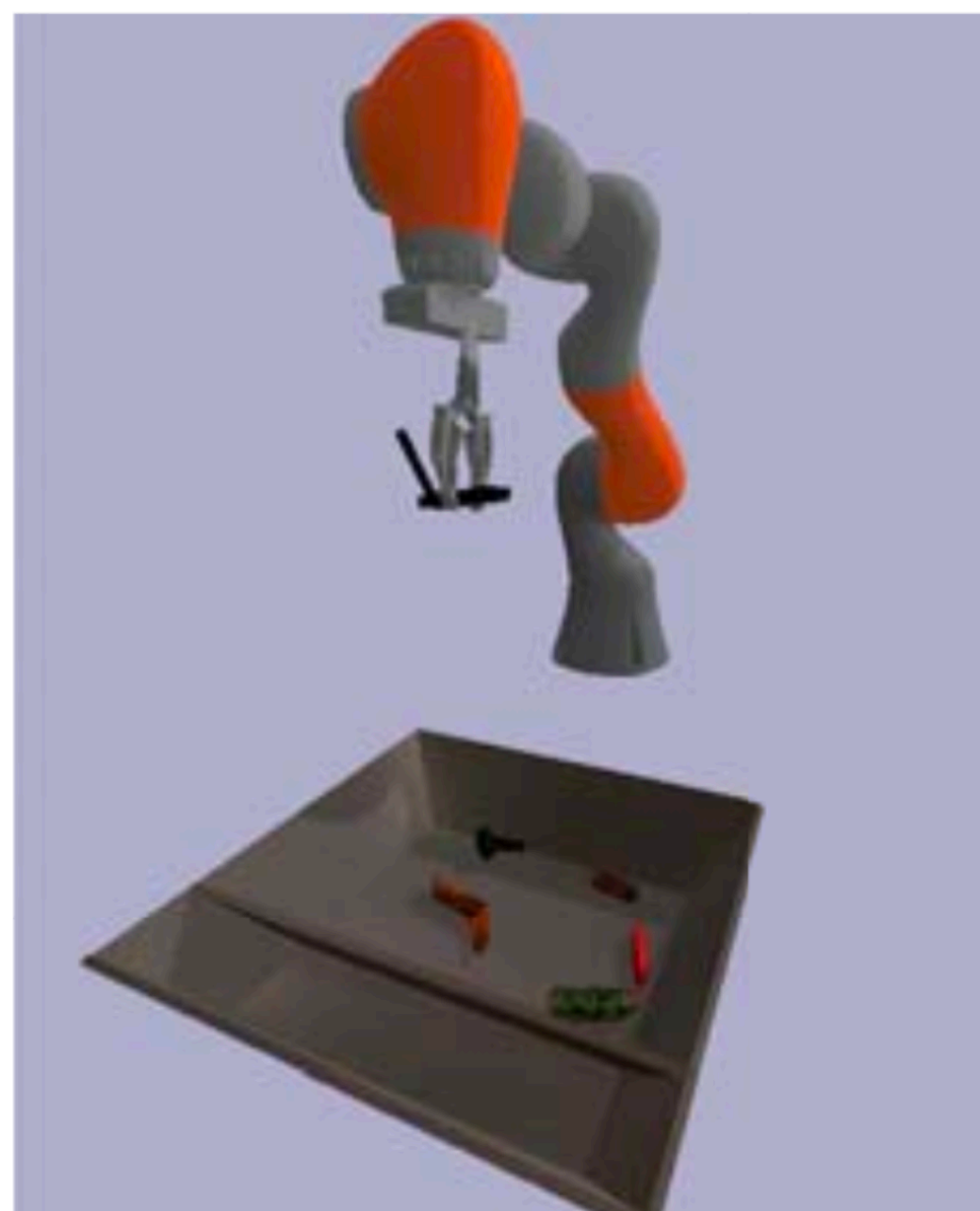
---

- Domain **randomization** alleviates the need for target **domain knowledge**
  - But much is still needed: a **simulator** in the ballpark, the randomization **ranges**
- The more we know about **target domain**, the better we can adapt the **source**
- Can we automate this **adaptation process**?
  - Using target-domain data

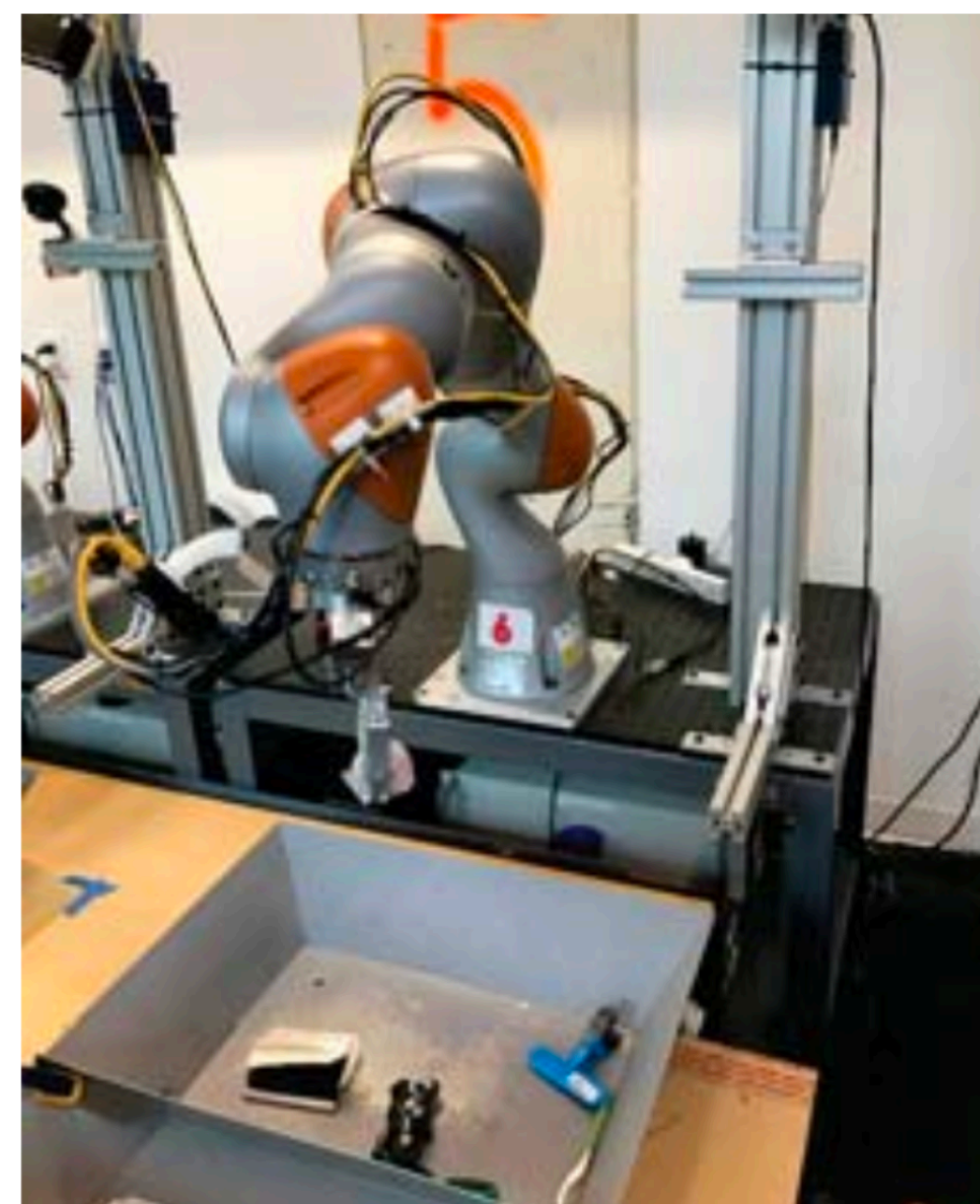


# Sim2real with domain adaptation

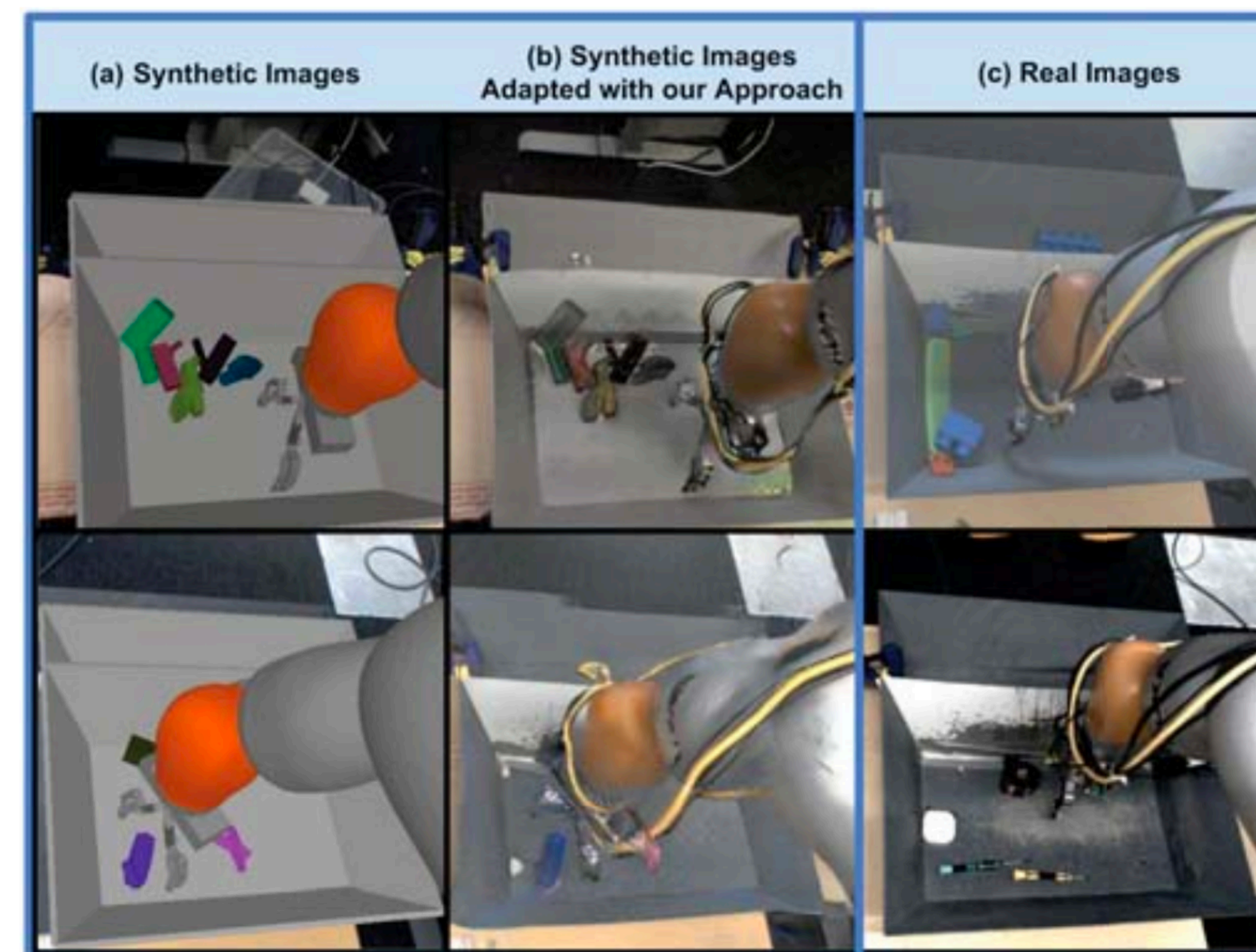
- Source domain looks **vaguely** like target domain  $\implies$  may not transfer well
- Idea: **adapt** the source domain to look more like the target (**more realistic**)



(a) Simulated World



(b) Real World





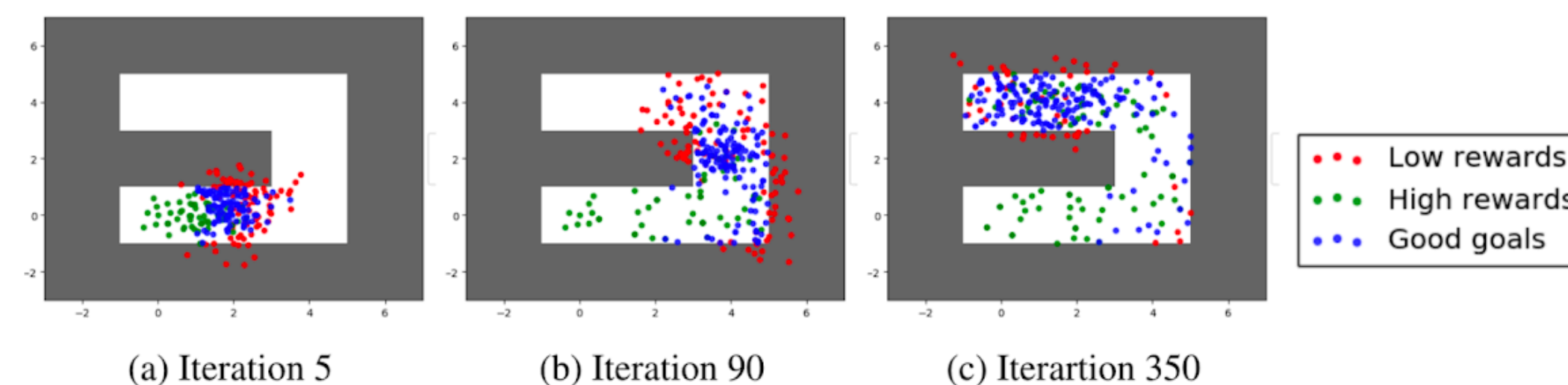
# Curriculum learning

---

- Why **pre-train** a policy?
  - ▶ So far: to collect **more data** in faster simulator
  - ▶ **Curriculum learning** = start with easy version of the task, make it **gradually harder**
- If a task is hard to solve by itself, “**training wheels**” can help
  - ▶ **Exploration** never finds rewards? **Shorten** task
  - ▶ Rewards don't encourage **exploring** / reaching **subgoals**? Leave “**breadcrumbs**”
  - ▶ Poor SGD **convergence** properties? **Coarsen** states / actions / time
  - ▶ Challenging state inference under **partial observability**? Add observability

# Goal GAN

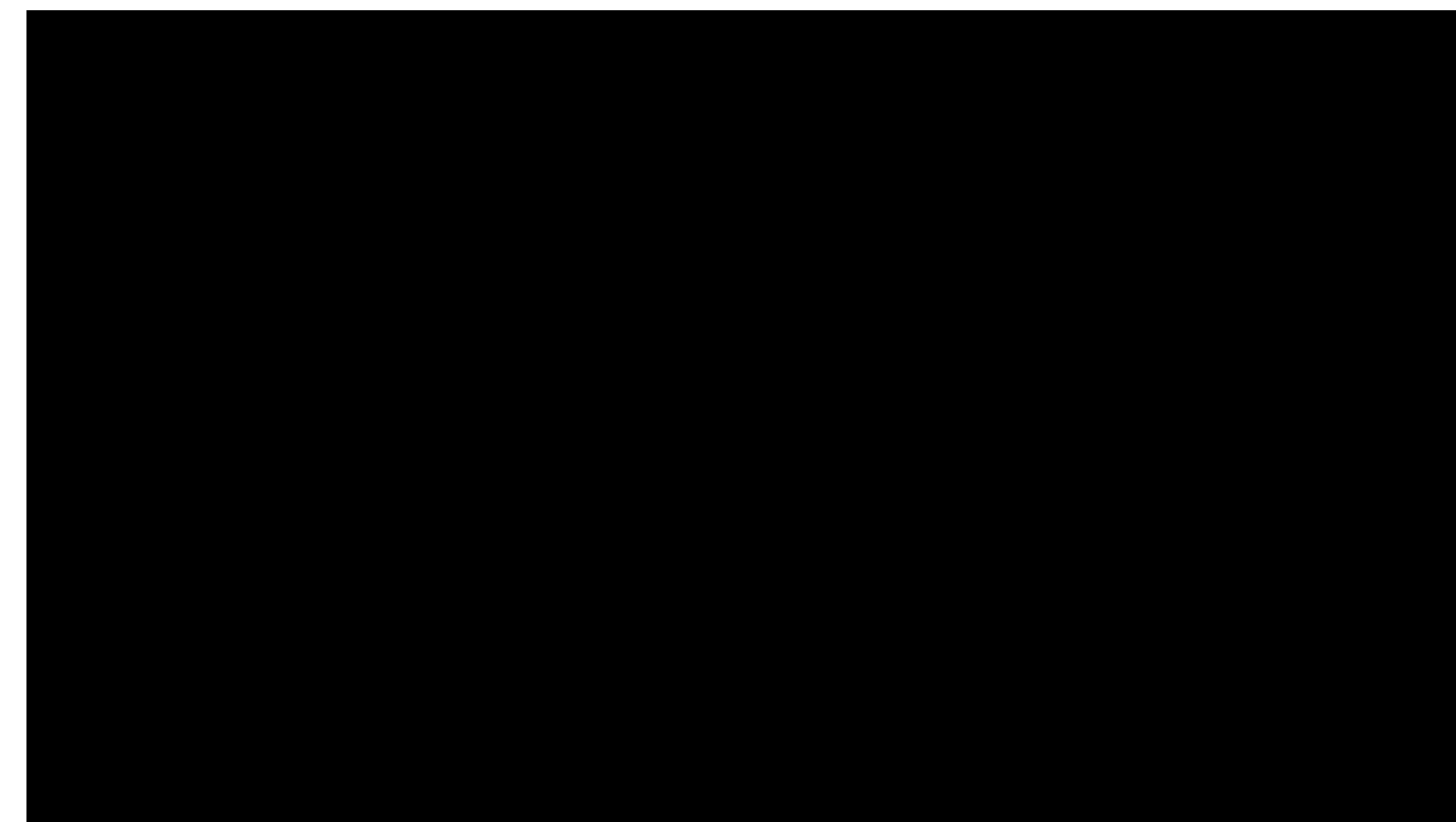
- **Sample** goals, roll out policy
  - ▶ **Reject** goals with too **low** / **high** rewards



- **Train GAN** to generate goals with this distribution
- **Train agent** on generated goals

- Generated goals are **intermediate-level**:

- ▶ just hard enough for the agent to **learn** something new
- ▶ not so hard that it **struggles** to do it



# Multi-task learning settings

---

- Transfer learning
  - Earlier domains / tasks are only stepping stones towards the ultimate task
- Shared learning
  - Learn multiple tasks jointly, have them inform each other
- Lifelong learning
  - Learn tasks as they occur, but also keep past abilities
    - Avoid **catastrophic forgetting** = fine-tuning a model **degrades** its quality for old task

# Today's lecture

---

Transfer learning

Domain randomization + adaptation

**Shared learning**

# Shared learning

---

- What can be **learned jointly** across tasks?
  - ▶ World **model** / perceptual **features**
    - Similar to above
  - ▶ **Policy**
    - Multi-task policy distillation — up next (1)
    - Task-aware policy — up next (2)
  - ▶ **Modules** in a structured policy
    - Multi-task hierarchical imitation learning (HIL-MT) — up next (3)



# Policy distillation

---

- **Policy distillation** = behavior cloning of existing policy with cross-entropy loss
- Wait, but... why?! if we already have the policy
  - ▶ Network **compression**
  - ▶ Track “**average**” policy (stabilize training, fictitious play in game theory, etc.)
  - ▶ **Combine** multiple policies into one — up next

# Multi-task policy distillation

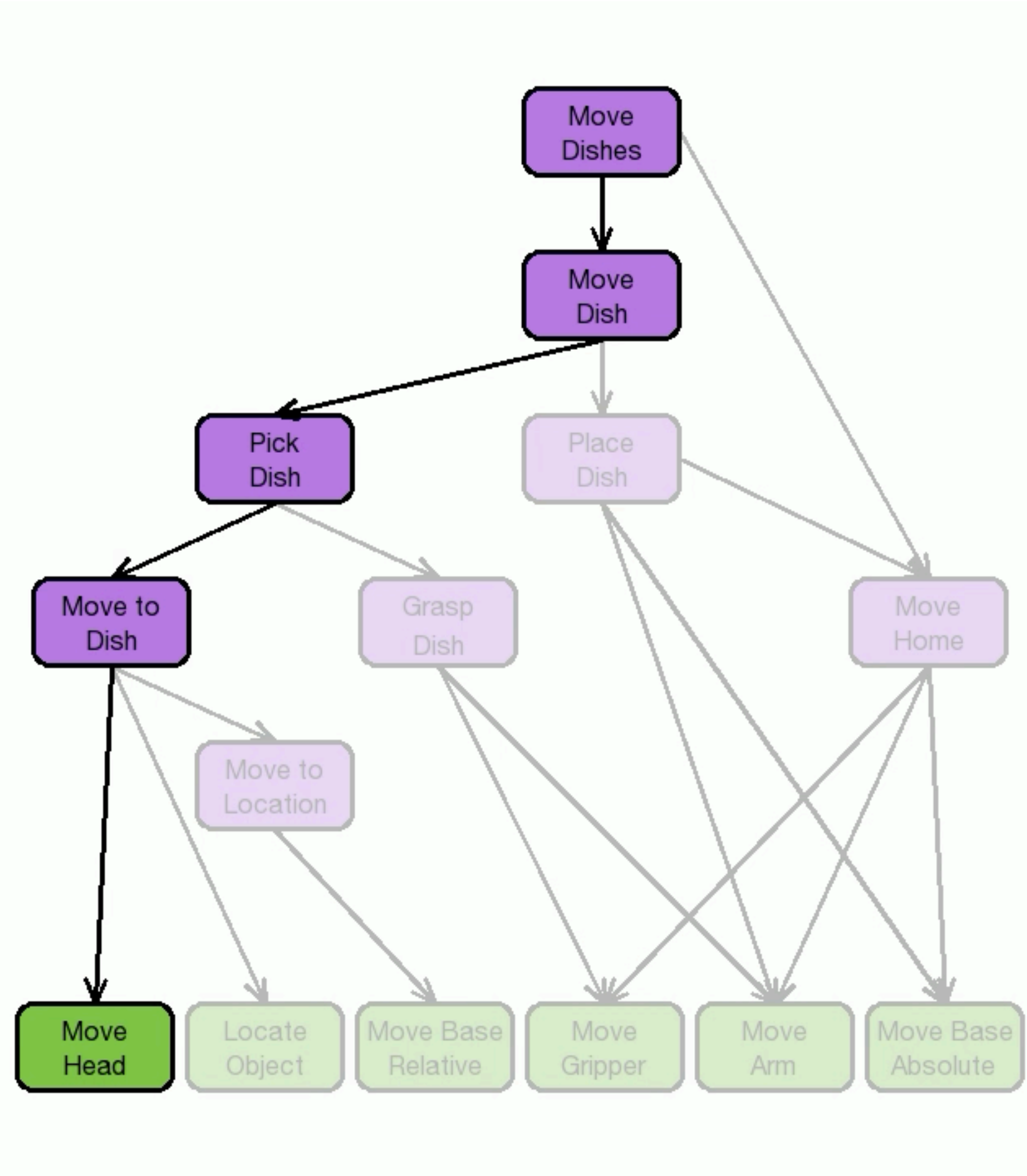
---

- **Train** a policy for each task
- **Distill** them into one policy
  
- **Pros:**
  - Effectively combining the empirical evidence from all datasets = **more data**
  - If tasks are related, distilled policy can be **more stable**
  
- **Cons:**
  - If tasks aren't related, they **compete** for network capacity
  - One very **wrong** distilled policy can ruin it for everyone

# Task-aware policy

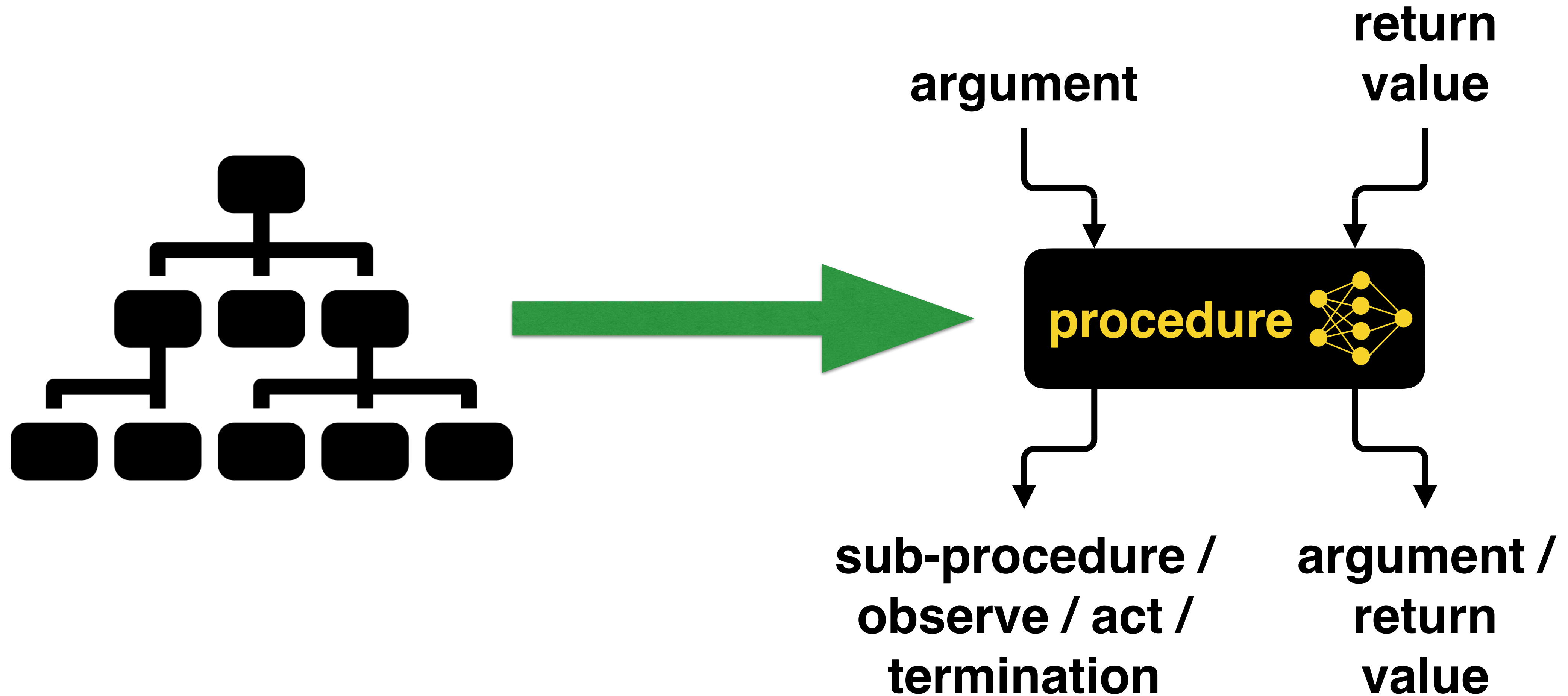
- Similar to **goal-conditioned** behavior cloning (Lecture 2)
  - Except: goal = state; task = (dynamics + reward)  $\implies$  more general
- Separate policy  $\pi_\tau(a | s)$  for each task  $\tau \rightarrow$  one task-aware policy  $\pi(a | s, \tau)$
- Sometimes, a task has a natural **embedding**
  - **Direction + speed** for Walker / Swimmer / Hopper
  - **Index** of a block to pick + **position** to place it
- Otherwise, can **learn to embed** task from specification
  - Task can be **specified** by demonstration / text / target image

# Learning from annotated demonstrations





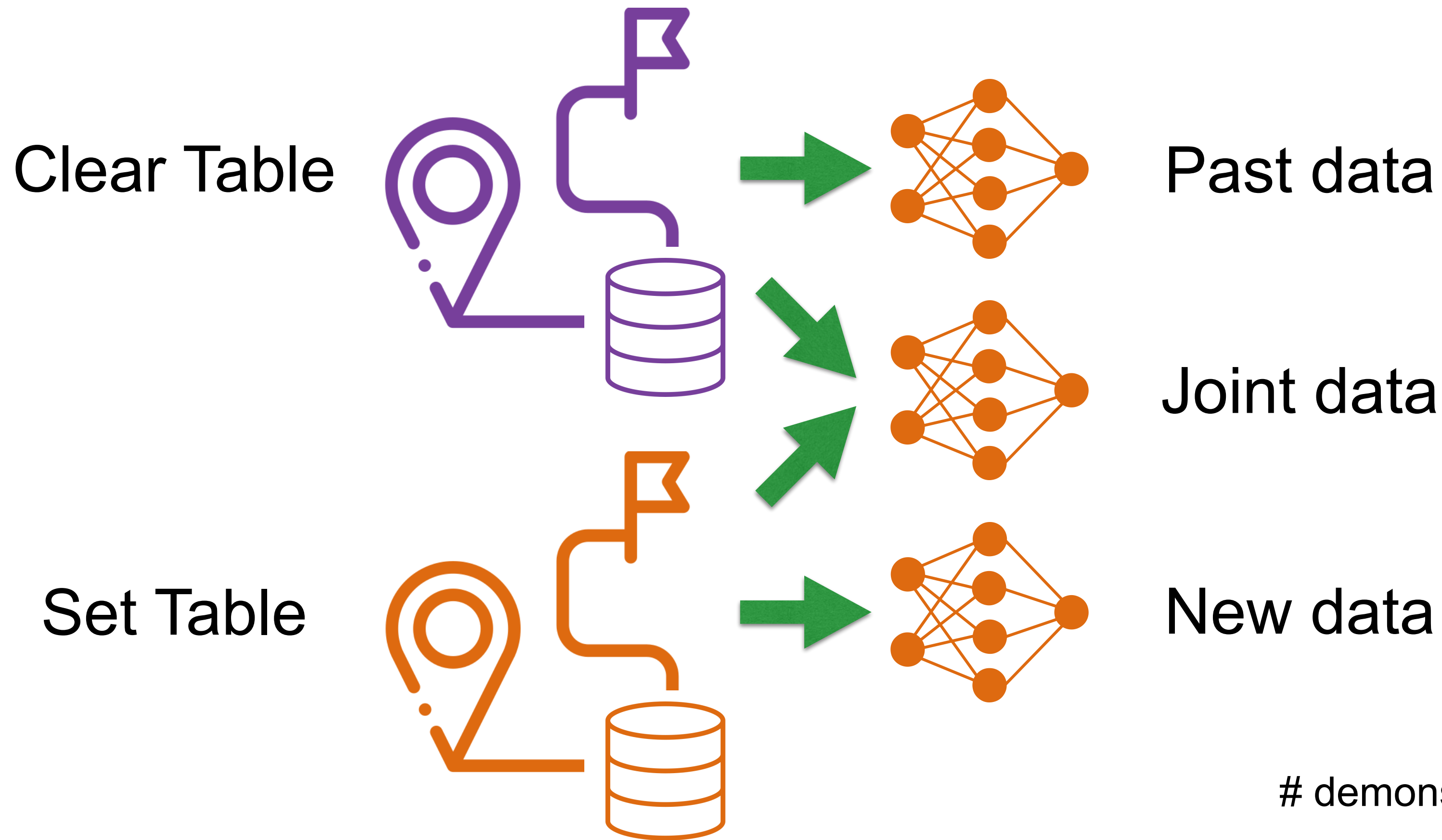
# Supervised learning each skill



$$\log p_{\theta}(\text{procedure steps}) = \sum_i \log p_{\theta}(\text{procedure step } i)$$

# Multi-task hierarchical imitation learning (HIL-MT)

- Hierarchical control allows per-procedure selection of multi-task mode



# demonstrations to learn Clear Table → Set Table

$\mathcal{D}_{clear}$	$\mathcal{D}_{set}$	$\mathcal{D}_{clear} \cup \mathcal{D}_{set}$	Per-skill selection
Failed	$19 \pm 0.3$	Failed	<b><math>11.6 \pm 0.25</math></b>

# Recap

---

- Reuse data between related tasks
  - May hurt if tasks are unrelated
- To improve the task overlap: soft-optimality, randomization, adaptation
- Shared learning may benefit both source and target tasks
- Modularity allows mix-and-match of best approach
- Did not talk about: meta-learning, lifelong learning