

# 1709.04905 - One-Shot Visual Imitation Learning via Meta-Learning

- Yunqiu Xu
  - Other reference:
    - Presentation: [https://www.youtube.com/watch?v=\\_9Ny2ghjwuY](https://www.youtube.com/watch?v=_9Ny2ghjwuY)
    - Code: <https://github.com/tianheyu927/mil>
- 

## 1. Introduction

- Challenge: learning each skill from a scratch is infeasible
- One-Shot Visual Imitation Learning via Meta-Learning
  - Reuse past experience and learn new skills from a single demonstration
  - Visual: use raw visual inputs
  - Meta-Learning: MAML [C. Finn et.al. 2017](#)
- Do not take task identity / demonstration as the input into a contextual policy: learn parameterized policy, adapt into new tasks through gradient updates

## 2. Related work

- In this work, the state of env is unknown → should be learned from raw sensory inputs
- Two challenges for learning from demonstrations:
  - Compounding errors: not in this work
  - **The need of a large number of demonstrations for each task**
- Inverse RL can reduce the number of demonstrations, but requires additional robot experience to optimize the reward [C. Finn et.al. 2016](#)
- How we tackle this: share data across tasks
  - First, use a dataset of demonstrations of many other tasks for meta learning,

then we can learn a new task from its single demonstration

- Some existed work:
  - Contextual policies: provide the task as an input to the policy or value function
  - Train a variety of controllers, then learn a mapping from given task to controller parameters
- What we use: **meta-learning**

### 3. Meta-Imitation Learning

Three requirements for what qualifies as a meta-learning system [2]:

- ❑ The system must include a learning sub-system, which adapts with experience
  - ❑ Experience is gained by exploiting meta knowledge extracted
    - ❑ In a previous learning episode on a single dataset
    - ❑ From different domains or problems
- By training for adaptation across tasks, meta-learning effectively treats entire tasks as datapoints.
  - MAML

# Model-Agnostic Meta-Learning

Learn the weights  $\Theta$  of a model such that gradient descent can make rapid progress on new tasks.

## Algorithm 1 Model-Agnostic Meta-Learning

**Require:**  $p(\mathcal{T})$ : distribution over tasks  
**Require:**  $\alpha, \beta$ : step size hyperparameters

- 1: randomly initialize  $\theta$
- 2: **while** not done **do**
- 3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$
- 4:   **for all**  $\mathcal{T}_i$  **do**
- 5:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  with respect to  $K$  examples
- 6:     Compute adapted parameters with gradient descent:  $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
- 7:   **end for**
- 8:   Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$
- 9: **end while**

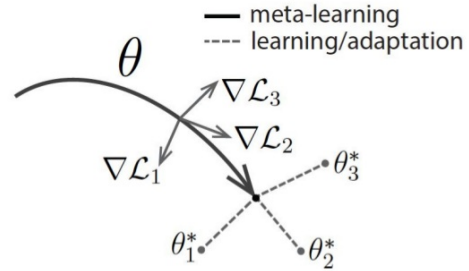


Figure 1. Diagram of our model-agnostic meta-learning algorithm (MAML), which optimizes for a representation  $\theta$  that can quickly adapt to new tasks.

$$\mathcal{T}_i = \{\tau = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_T, \mathbf{a}_T\} \sim \pi_i^*, \mathcal{L}(\mathbf{a}_{1:T}, \hat{\mathbf{a}}_{1:T}), T\}$$

Experts

- Extend MAML to MIL

Learn a policy that can quickly adapt to new tasks from a single demonstration of that task.

## Algorithm 1 Model-Agnostic Meta-Learning

**Require:**  $p(\mathcal{T})$ : distribution over tasks  
**Require:**  $\alpha, \beta$ : step size hyperparameters

- 1: randomly initialize  $\theta$
- 2: **while** not done **do**
- 3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$
- 4:   **for all**  $\mathcal{T}_i$  **do**
- 5:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  with respect to  $K$  examples
- 6:     Compute adapted parameters with gradient descent:  $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
- 7:   **end for**
- 8:   Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$
- 9: **end while**

## Algorithm 1 Meta-Imitation Learning with MAML

**Require:**  $p(\mathcal{T})$ : distribution over tasks  
**Require:**  $\alpha, \beta$ : step size hyperparameters

- 1: randomly initialize  $\theta$
- 2: **while** not done **do**
- 3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$
- 4:   **for all**  $\mathcal{T}_i$  **do**
- 5:     Sample demonstration  $\tau = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_T, \mathbf{a}_T\}$  from  $\mathcal{T}_i$
- 6:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  using  $\tau$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation (2)
- 7:     Compute adapted parameters with gradient descent:  $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
- 8:     Sample demonstration  $\tau'_i = \{\mathbf{o}'_1, \mathbf{a}'_1, \dots, \mathbf{o}'_T, \mathbf{a}'_T\}$  from  $\mathcal{T}_i$  for the meta-update
- 9:   **end for**
- 10:   Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$  using each  $\tau'_i$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation 2
- 11: **end while**
- 12: **return** parameters  $\theta$  that can be quickly adapted to new tasks through imitation.

$$\mathcal{T}_i = \{\tau = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_T, \mathbf{a}_T\} \sim \pi_i^*, \mathcal{L}(\mathbf{a}_{1:T}, \hat{\mathbf{a}}_{1:T}), T\}$$

Experts

- Loss function (Equation 2 in pseudo code)

$$\mathcal{L}_{\mathcal{T}_i}(f_{\phi}) = \sum_{\tau^{(j)} \sim \mathcal{T}_i} \sum_t \|f_{\phi}(\mathbf{o}_t^{(j)}) - \mathbf{a}_t^{(j)}\|_2^2. \quad (2)$$

- Network architecture

**Loss function:**

$$\mathcal{L}_{\mathcal{T}_i}(f_\phi) = \sum_{\tau^{(j)} \sim \mathcal{T}_i} \sum_t \|f_\phi(\mathbf{o}_t^{(j)}) - \mathbf{a}_t^{(j)}\|_2^2.$$

**Hyperparameters:**

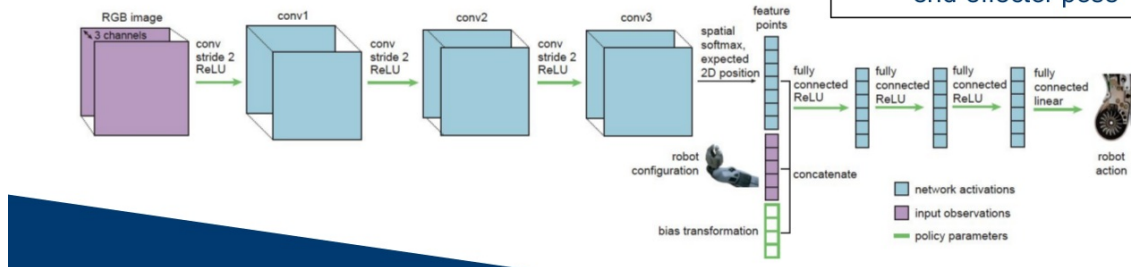
Conv: 40 3x3 filters

Fully connected: 200 dimension

Robots configuration:

joint angles,

end-effector pose



- Training and testing

## Meta-Imitation Learning

*[meta-training time]* 
$$\min_{\theta} \sum_{\text{tasks}} \mathcal{L}_v(\theta - \alpha \nabla_{\theta} \mathcal{L}_{tr}(\theta))$$

meta-training tasks

val demo

training demo

*[meta-test time]* 
$$\theta' \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta)$$

demo of meta-test task

**behavioral cloning loss**

$$\mathcal{L} = \sum_t \|\pi_{\theta}(\mathbf{o}_t) - \mathbf{a}_t^*\|^2$$

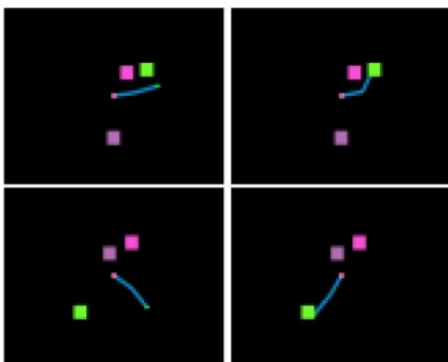
37

## 4. Experiment

## Comparison:

- ❑ **Random policy:** A policy that outputs random actions from a standard Normal distribution
- ❑ **Contextual policy:** A feedforward policy, which takes as input the final image of the demonstration and the current image, and outputs the current action.
- ❑ **LSTM:** use a recurrent network to ingest the provided demonstration and current observation, and outputs the current action.
- ❑ **LSTM+attention:** use attention mechanism

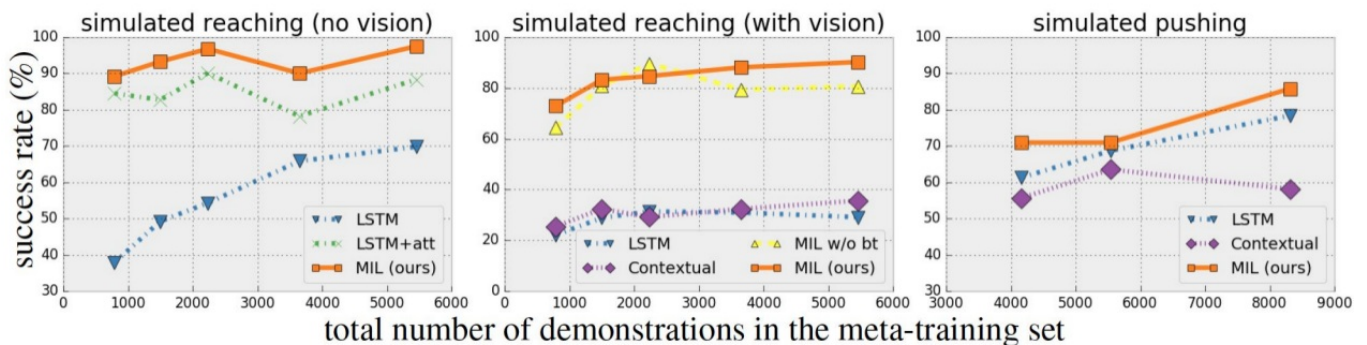
## Simulated reaching



### Task:

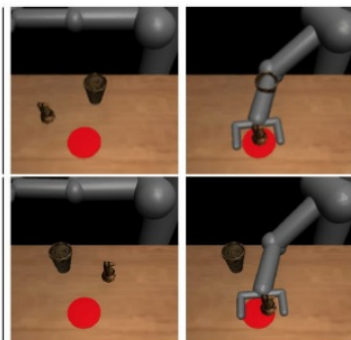
- ❑ reaching a target of a particular color
- ❑ Policy must learn to localize the target using the demonstration and generalize to new positions
- ❑ Meta-training must learn to handle different colors
- ❑ 150 tasks x 10 different trials per task

### Simulated reaching





# Simulated pushing

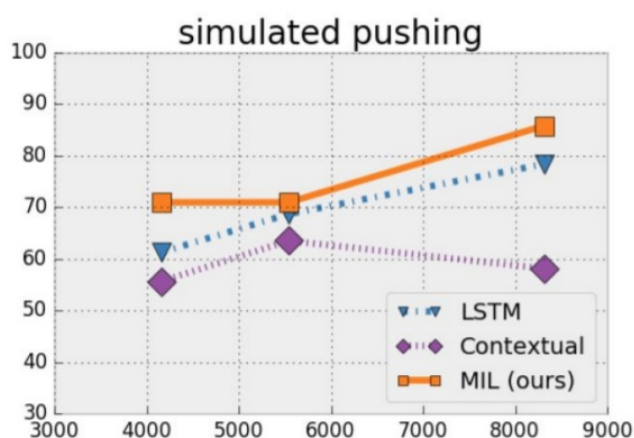


## Task:

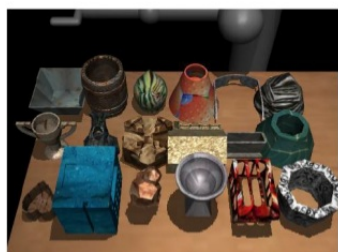
- A push is considered as success if the center of the target object lands on the red target circle for at least 10 timestamps.
- Each task is defined as pushing a particular objects
- 74 tasks x 6 different trials per task

method		video+state +action	video +state	video
LSTM	1-shot	78.38%	37.61%	34.23%
contextual		n/a	58.11%	56.98%
MIL (ours)		<b>85.81%</b>	<b>72.52%</b>	<b>66.44%</b>
LSTM	5-shot	83.11%	39.64%	31.98%
contextual		n/a	64.64%	59.01%
MIL (ours)		<b>88.75%</b>	<b>78.15%</b>	<b>70.50%</b>

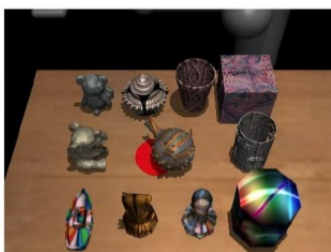
Table 1: One-shot and 5-shot simulating pushing success rate with varying demonstration information provided at test-time. MIL can more successfully learn from a demonstration without actions and without robot state and actions than LSTM and contextual policies.



# Real-World Placing



subset of training objects



test objects



subset of training objects



test objects

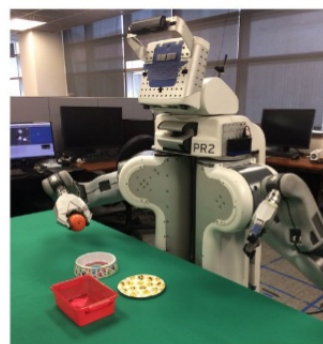
## Task:

Evaluate how well a real robot (PR2) can learn to interact with new unknown objects from a single visual demonstration.

**Success:** the held object landed in or on the target container after the gripper is opened

method	test performance
LSTM	25%
contextual	25%
MIL	<b>90%</b>
MIL, video only	<b>68.33%</b>

Table 2: One-shot success rate of placing a held item into the correct container, with a real PR2 robot, using 29 held-out test objects. Meta-training used a dataset with ~100 objects. MIL, using video only receives the only video part of the demonstration and not the arm trajectory or actions.



- Demo: let robot put the ball in blue cup
- Test: shuffle the cups and let robot achieve the goal

## 5. Summary and Ongoing Work (On CoRL)

- Summary:
  - reuse prior experience when learning in new settings
  - learning-to-learn enables effective one-shot learning
- Ongoing work: one-shot imitation from human video → during demo, let human put the ball in cup