

---

# Characterizing Policy Divergence for Personalized Meta-Reinforcement Learning

---

Michael Zhang

Harvard University

Cambridge, MA 02138

michael\_zhang@college.harvard.edu

## Abstract

Despite ample motivation from costly exploration and limited trajectory data, rapidly adapting to new environments with few-shot reinforcement learning (RL) can remain a challenging task, especially with respect to personalized settings. Here, we consider the problem of recommending optimal policies to a set of multiple entities each with potentially different characteristics, such that individual entities may parameterize distinct environments with unique transition dynamics. Inspired by existing literature in meta-learning, we extend previous work by focusing on the notion that certain environments are more similar to each other than others in personalized settings, and propose a model-free meta-learning algorithm that prioritizes past experiences by relevance during gradient-based adaptation. Our algorithm involves characterizing past policy divergence through methods in inverse reinforcement learning, and we illustrate how such metrics are able to effectively distinguish past policy parameters by the environment they were deployed in, leading to more effective fast adaptation during test time. To study personalization more effectively we introduce a navigation testbed to specifically incorporate environment diversity across training episodes, and demonstrate that our approach outperforms meta-learning alternatives with respect to few-shot reinforcement learning in personalized settings.

## 1 Introduction

While reinforcement learning (RL) has demonstrated success in various sequential decision making settings Barto et al. [1990], Shortreed et al. [2011], challenges remain in the personalized domain. For example, we may be interested in recommending treatments to a cohort of patients, where individuals may respond to the same treatment differently, even while conditioning on the same health state. From an RL framework, each patient can be modeled as an entity whose individual characteristics parameterize a distinct Markov decision process (MDP). With this framing, one could, in principle, train a separate policy from scratch for each individual. However, this approach is impractical in settings where exploration is costly and individual trajectory data are limited. As a result, there is an increasing need to think about how to leverage prior data in sophisticated manners to develop few-shot learning approaches for deployment in unseen environments.

In this work, we take a meta-learning approach to this problem and propose a method tailored to few-shot learning across a *diverse* set of environments. Our approach builds on model agnostic meta-learning (MAML), a paradigm for learning a single prior initialization of policy gradient parameters that can adapt quickly to a set of observed tasks. We extend this approach by leveraging the notion that certain environments are more similar to each other than others in personalized settings, and accordingly their corresponding experiences should be given higher attention during adaptation. We

thus propose generalizing this approach by learning several different potential initializations, and choosing an appropriate initialization to use for a given environment at test time.

Our proposed algorithm—cluster adaptive meta-learning (CAML)—explicitly seeks to determine which past parameters should be factored in during adaptation. Because the core idea of our adaptation relies on being able to quickly identify the most relevant past trajectories, we seek to characterize these experiences through a distance metric that establishes a notion of similarity between policies, borrowing from related work in inverse reinforcement learning (IRL) and determining policy divergence. We rely on whole trajectories—which provide a richer source of information over alternatives used in previous work Clavera et al. [2018], Finn et al. [2017], Nagabandi et al. [2018]—and would like to call upon a diverse set of environment interactions for reference. Accordingly, we retain an episodic model-free approach in contrast to more recent model-based online-learning methods where each timestep is considered a new task Clavera et al. [2018]. Our contributions center around how to (1) study personalization effectively in simulated contexts, (2) characterize the divergence of policies across a population of environments, and (3) learn effectively across different and previously unseen environments with minimum exploration cost. We introduce a  $K$ -medoids-inspired algorithm for policy adaptation and few-shot learning across multiple environments, and demonstrate approach’s successes over existing methods in a 2D testbed.

## 2 Background

### 2.1 Personalized Markov decision processes (MDPs)

We consider a modified version of a typical continuous state-space MDP, parameterized specifically by an individual entity type  $\mathcal{T}_i$  among a population of possible types  $\mathcal{T}$ . For lexical consistency, we refer to the agent as the actual policy or decision-maker, and define an “entity” to be the object through which an agent interacts with the larger environment (the “patient” in our working example). Each entity type then introduces an MDP  $\mathcal{M}^i \triangleq (\mathcal{S}, \mathcal{A}, P^i, r^i, \gamma, s_0^i)$  with unique transition probabilities  $P^i : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ . Given shared state and action spaces  $\mathcal{S}$  and  $\mathcal{A}$ , we thus allow for diversity in behavior across entities from different types.

For any two different entity types  $\mathcal{T}_i \neq \mathcal{T}_j$ ,  $P^i$  may not necessarily differ from  $P^j$ , but crucially at the onset of training for some subset of unseen types  $\mathcal{T}_i \subseteq \mathcal{T}$  we make no assumption that these transition dynamics are the same. Finally,  $\gamma$  is a discount factor assumed to be the same for all agents, and  $s_0^i$  is the initial state distribution for environment type  $i$ . Our goal across a population of entity types is then to learn some function  $f : \mathcal{T} \rightarrow \Pi$  mapping from entity types to policies. For each entity type given  $s_0^i$  and a minimal number of timesteps  $t$ , we can then output an optimal policy  $\pi_\theta^i$  with regard to maximizing the cumulative discounted reward  $\sum_{t=0}^{T-1} \gamma^t r(s_t, a_t)$  for  $T$ -length episodes, where  $\theta$  denotes the parameters for personalized policy  $\pi_i$  taking action  $a_t \in \mathcal{A}$  given state  $s_t \in \mathcal{S}$ .

### 2.2 Meta-reinforcement learning

In the typical meta-learning scenario, we are interested in automatically learning learning algorithms that are more efficient and effective than learning from scratch Clavera et al. [2018]. To do so, instead of looking at individual data points, a meta-learning model is trained on a set of tasks, treating entire tasks as training examples with the goal being to quickly adapt to new tasks using only a small number of examples at test time. We can consider a task to be any type of desirable machine learning activity, such as correctly classifying cat images in a supervised learning scenario, or learning how to walk forwards at a certain velocity in RL. Expanding on this latter example, meta-RL tries to learn a policy for new tasks using only a small amount of experience in the test setting. More formally, each RL task  $\tau_i$  contains an initial state distribution  $s_0^i$  and transition distribution  $P^i(s_{t+1}|s_t, a_t)$ . In the few-shot RL setting parameterized by  $K$  shots, after training on train tasks  $\{\tau_{\text{train}}\}$ , we are allowed to use  $K$  rollouts on test task  $\tau_{\text{test}}$ , trajectory  $(s_0, a_0, \dots, s_T)$ , and rewards  $R(s_t, a_t)$  for adapting to  $\tau_{\text{test}}$ .

### 2.3 Model-agnostic meta-learning for RL

More specifically, our approach builds on model-agnostic meta-learning (MAML) for reinforcement learning Finn et al. [2017], which optimizes over the average reward of multiple rollouts. Doing so we wish to learn an initial set of policy gradient parameters capable of quickly adapting to a

new RL task in a few gradient descent steps. Given  $n \in N$  meta-training environments, and initial policy parameters  $\theta$ , MAML first samples an MDP  $\mathcal{M}^i$ ,  $i \in \{1, 2, \dots, n\}$ . For each sampled MDP, we then train a policy given  $K$  rollouts, arriving at task-adapted parameters  $\theta'_i$  through gradient descent  $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(\pi_{\theta})$  (described as the inner training loop).  $\alpha$  is our learning rate and  $\mathcal{L}_{\mathcal{T}_i}(\pi)$  describes the inverse reward from following policy  $\pi$  with entity type  $\mathcal{T}_i$ . Feedback through the reward generated on the  $(K + 1)^{\text{st}}$  rollout is then saved. After  $n$  tasks are sampled, we perform a meta-update using each saved reward, updating our initialization parameters  $\theta'$  with  $\theta' \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(\pi_{\theta'_i})$ , where  $\beta$  is the meta-learning rate. MAML thus seeks to leverage the gradient updates for each task during the inner loop of training, reaching a point in parameter space relatively close to the optimal parameters of all potential tasks. We then evaluate performance by sampling  $\mathcal{M}^j$ ,  $j \in N \setminus n$ , observing the rewards generated from the  $(K + 1)^{\text{st}}$  trajectory after adapting with  $K$  rollouts and policy parameters  $\theta'$ .

## 2.4 Characterizing policy divergence

Aside from gaining further knowledge with respect to how policies adapt to their environments over time, we also seek to establish a valid distance metric between various policies to inform policy parameter initialization in few-shot RL. Conceptually we can imagine that given pairs of MDPs and policies  $(\mathcal{M}^i, \pi^i)$  for types  $\mathcal{T}_i \in \mathcal{T}$ , the policies start with an indistinguishable parameter initialization and diverge when adapting to their personalized environments after  $t$  timesteps of training. If adapted policies  $\pi'_i$  and  $\pi'_j$  are still similar to each other, then this suggests that their environments  $\mathcal{M}^i$  and  $\mathcal{M}^j$  may be similar as well, and accordingly we can use this information to form a representation of optimal policy parameter clusters, making sure to only reference relevant prior experiences when adapting to a new environment.

To begin our search for such a distance metric we look to the field of imitation learning, and specifically inverse reinforcement learning (IRL). Here the goal is to learn a cost function explaining an expert policy’s behavior Ng and Russell [2000]. However, while typical IRL methods first learn this cost function and then optimize against it to train a policy to imitate the expert, we are content with the cost function alone, whose components we can borrow to define our own quantifiable metric defining the distance between two policies. As noted in Ho and Ermon [2016], we can uniquely characterize a policy by its occupancy measure given by

$$\rho_{\pi}(s, a) = \pi(a|s) \sum_{t=1}^T P(S_t = s|\pi), \text{ where } \rho_{\pi} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} \quad (1)$$

We essentially view this metric as the state-action distribution for a policy. Using a symmetric measure such as the Jensen-Shannon (JS) divergence Lin [1991], we can thus compare the occupancy measures of multiple policies at specific time steps over training as a metric into policy divergence over time. While we study a very different problem from traditional IRL—where we do not explicitly try to minimize policy divergence—such measures of divergence allow for both descriptive analysis in tracking divergence and informing new learning algorithms. In principle we can then characterize policies based on their observed trajectories, and later show that this is enough to measure similarities across their respective environments as well.

## 3 Related Work

While there is much work summarizing the general meta-learning literature Clavera et al. [2018], Hsu et al. [2018], Ren et al. [2018], we focus on initialization-based methods for RL. Here we wish to learn a set of optimal model parameters such that after meta-learning, the policy is initialized to an optimal position in parameter space to adapt to new environments relatively quickly. Finn et al. demonstrate this paradigm in MAML, which achieves meta-reinforcement learning through gradient-based optimization Finn et al. [2017]. They introduce the generalizable notion of doing well across a variety of tasks, where each task in the RL setting is a similar goal such as trying to reach a certain point in 2D space. Given training on reaching a certain support set of points, we would like to be able to quickly learn how to reach a new batch of unseen query points. Instead of trying to generalize across strictly different tasks, we adopt this framework by modeling different entity types encountered—each harboring their own potentially unique rewards and transition probabilities given the same states and actions—as different tasks.

Implementation-wise, although modern day software packages make MAML’s computation relatively straightforward, it’s reliance on unrolling the computation graph behind gradient descent and taking second derivatives has motivated follow-on work on simpler methods, such as first-order MAML and Reptile, which avoid doing so with first-order approximations and direct gradient movements respectively Finn et al. [2017], Nichol et al. [2018]. We focus on learning algorithms with similar computational dependencies due to their comparable performance with vanilla MAML.

Learning similarities between tasks through clustering or other distance-based metrics has also been explored before in non-RL few-shot settings, where a central motivation lies in being able to achieve comparable results with alternatives while maintaining much simpler inductive biases. Snell et al. propose prototypical networks for few-shot classification Snell et al. [2017], which aims to learn an embedding where data points cluster around a single prototype representation in each class. Accordingly, given a learned transformation between input features and some vector space, we can simply use the computed prototypes on unseen data points for classification. Although we do not explore learning further embeddings on our proposed distance metric, we nonetheless demonstrate that calculating similarities given observed occupancy measures is enough to identify similar environments, drawing upon a similar intuition for RL.

## 4 Personalizing policy learning

We now present our approach for (1) measuring and interpreting divergence of policies in personalization, and (2) defining few-shot meta-learning algorithms for online adaptation to new entity types. We also introduce a testbed to study the effects of personalization across entities on training. As outlined in Section 2.1, we work with a set of personalized MDPs representing a population of entity types.

### 4.1 A personalized particle environment for 2D navigation

We begin our study of divergence in personalizing policies with a 2D continuous episodic gridworld environment, where a point agent must move to a target position in a constrained amount of time. The state space is given by the agent’s current 2D position, and at each timestep the agent may choose to go left, right, up, or down one pixel unit. Following these commands, the agent tries to maximize its reward, calculated to be the negative squared distance to the goal after each timestep. Finally episodes end every 100 timesteps, or when the agent is within 0.01 units of the goal. Particles start at origin  $(0, 0)$ , and try to reach target at  $(1, 1)$  before the episode ends.  $\mathcal{S}$  is defined by the coordinates of the particle at timestep  $t \in \{0, 1, \dots T\}$ . Actions are discrete  $a$ , corresponding to moving one unit left, right, down, up, i.e.  $a_t \in \{(-1, 0), (1, 0), (0, -1), (0, 1)\}$ .

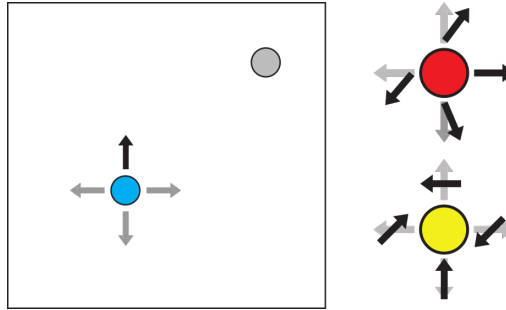


Figure 1: **2D personalized particles.** Policy  $\pi$  tries to move entity (blue, red, yellow) to target (gray), but entities behave in different ways unknown to the policy. Right: Two remapping schemes, where some cardinal direction e.g. ‘up’ corresponds to a different transition vector. Code available upon request.

We personalize this environment with the introduction of a population of entities, where we may either stick with one throughout the policy’s entire training run or randomly introduce new entities into the world at the start of every episode. Each agent is initialized with a personalization function  $\mathcal{F} : \mathcal{A} \rightarrow \mathcal{A}'$ , which first remaps the cardinal directions of each action and then imposes additional variance. For example, given default action  $a_r = (1, 0)$ , telling a default agent to go right, our policy may encounter an agent who upon receiving action input  $(1, 0)$  actually moves with action  $(0.42, -1)$ . This is achieved by first remapping  $a_r$  to  $(0, -1)$ , and then imposing further variance on the x-coordinate. We can think of this as having the same set of interventions for all agents, but facing different responses in the form of varied actual transition outcomes.

## 4.2 Practical divergence estimation

As introduced in Section 2, by characterizing policies through their occupancy measures, we derive a similarity metric and quantitatively analyze their divergence. Furthermore, individual policies may move closer or further to others over time, forming new neighbors that reflect personalization over time. However, as discussed in [Ho and Ermon, 2016], simply using the occupancy measure as presented and matching on all states and actions is not practically useful. Instead of looking at an occupancy measure defined over all states as above then, we propose an alternative measure reliant only on a sampled subset of states. We observe that two policies may intuitively be different if for each of various given states, they differ in their action distributions. However, the states must have a non-zero chance of occurring in both policies of interest. Accordingly, in the  $K$ -shot  $N$ -ways setting, where for each of  $N$  MDPs  $\mathcal{M}_i$  we get  $K$  trajectories  $\{\tau_i^1, \dots, \tau_i^K\}$  to observe before updating and an additional trajectory  $\tau_i'$  after the update, we can calculate a kernel density estimate (KDE)  $\hat{q}(\tau_{1:N})$  based on the observed trajectories  $\tau_{1:N}' = \{\tau_1', \tau_2', \dots, \tau_N'\}$ . States  $\hat{s}_j \sim \hat{q}(\tau_{1:N}')$  can then be sampled, with probability density  $\hat{q}(\tau_{1:N}')(\hat{s}_j)$ .

Referring back to (1), for any policy  $\pi_i$ , we estimate the conditional probability  $\pi(\cdot|s)$  by feeding in  $\hat{s}_j$  to our model network and obtaining the probability vector  $\pi_i'(\hat{s}_j) = (p_1^j, \dots, p_n^j)$  where  $p_a^j$  denotes the probability of taking action  $a$  given state  $s_j$ . Towards estimating  $\sum_{t=1}^T P(S_t = s|\pi)$ , we again calculate the KDE over our observed trajectories, but now only consider  $\tau_i'$  for MDP  $\mathcal{M}_i$ .

Accordingly, for any set of policies we first obtain a respective observed batch of trajectories, calculate an estimated state-marginalized variant of occupancy measure  $\hat{\rho}_s^i$  for  $s \in \mathcal{S}$  and entity types  $\mathcal{T}_i$ , and compute pairwise symmetric divergences across entity types. In our case we use the Jensen-Shannon divergence Lin [1991].

$$D_{JS}(\rho_s^i, \rho_s^j) \triangleq D_{KL}(\rho_s^i \parallel \frac{\rho_s^i + \rho_s^j}{2}) + D_{KL}(\rho_s^j \parallel \frac{\rho_s^i + \rho_s^j}{2}) \text{ for } i \neq j \quad (2)$$

where  $D_{KL}(\cdot, \cdot)$  is the Kullback–Leibler divergence. As a final measure, we employ  $K$ -medoids clustering to group our policies using the pairwise  $D_{JS}(\cdot, \cdot)$  as a distance metric. The overall procedure is summarized in Algorithm 1. Besides describing our policies, we use these metrics to inform algorithms for few-shot learning, as described in the next section.

---

### Algorithm 1 Policy Divergence Estimation through Observed Trajectories

---

**Require:** Distribution over entity types  $p(\mathcal{T})$ , initial policies  $\pi_\theta^i \in \Pi$  and episode horizon  $T$

**Require:** Initialize update counter  $t = 0$ , and batch size  $n$  for computing estimates

```

1: while not done do
2:   Sample entity type  $\mathcal{T}_i \sim p(\mathcal{T})$  and initialize new policy  $\pi_\theta^i$ 
3:   Sample  $K$  trajectories  $\tau^i = \{s_1, a_1, \dots, s_T\}$  following  $\pi_\theta^i$ 
4:   Update  $\theta'$  through vanilla policy gradient with REINFORCE Sutton et al. [2000]
5:   Save trajectory  $\tau_i' = \{s_1, a_1, \dots, s_T\}$  using updated  $\pi_\theta^i$ 
6:   Increment  $t \leftarrow t + 1$ 
7:   if  $t = n$  then
8:     Compute type-specific KDE  $\hat{q}(\tau_i')$  for all saved  $\tau_i'$ 
9:     Obtain concatenated updated trajectories  $\tau_{1:N}' = \{\tau_1', \dots, \tau_N'\}$ 
10:    Compute overall KDE  $\hat{q}(\tau_{1:N})$  and generate samples  $s \sim \hat{q}(\tau_{1:N})$ 
11:    for all  $s$  do
12:      Calculate  $\pi^i(a|s)$  for all policies  $\pi^i$  and compute state probability densities  $\hat{q}(\tau_i')(s)$ 
13:      Estimate  $\hat{\rho}_s^i = \pi^i(a|s)\hat{q}(\tau_i')(s)$  for given state  $s$ 
14:    end for
15:    Compute time-indexed pairwise Jensen-Shannon divergences  $D_{JS}(\hat{\rho}_s^i, \hat{\rho}_s^j)$  for all  $\mathcal{T}_i, \mathcal{T}_j \in \mathcal{T}$ 
16:    Obtain divergence metric  $\sum_{s \sim \hat{q}(\tau_{1:N}')} D_{JS}(\hat{\rho}_s^i, \hat{\rho}_s^j)\hat{q}(\tau_{1:N}')(s)$ 
17:    Reset  $t = 0$ 
18:  end if
19: end while

```

---

### 4.3 Cluster-adapting meta-learning

Given our policy divergence estimators, we hope to better organize our previous experiences for adaptation in new environments. In this section, we describe a  $K$ -medoids-inspired meta-learning algorithm called cluster-adapting meta-learning (CAML). Similar to serial versions of meta-learning algorithms such as MAML, CAML iteratively learns an initialization for parameters of a neural network model, such that given new MDP  $\mathcal{M}^i$ , after a few trajectory rollouts and a single batch update during test time, our policy performs competitively to those pretrained on the same environment. However, while previous algorithms update a single set of optimal parameters drawn from parameter space, CAML maintains parameters representative of larger clusters.

In setting up CAML, we consider the following design choices: (1) a valid distance metric to compare past trajectories and their corresponding policies in a model-free setting, (2) the clustering algorithm to organize past experiences, and (3) an initialization method to select the most appropriate policy parameters at test time. We consider (1) to be the most important contribution, and included details in Section 4.2. As described, we implement (2) using the  $K$ -medoids algorithm Kaufmann [1987], although in practice any clustering method may be used. Finally with regard to (3) at the end of our training iterations we obtain  $k$  medoid policies, and during evaluation we view fast-adaptation to each individual MDP as a bandit problem. Given  $k$  available arms (the policies) and  $K$  arm pulls (the few number of shots, i.e. episodes allowed), we want to maximize the corresponding reward from following the arm policy for an entire episode. While more sophisticated multi-armed bandit methods exist Merentitis et al. [2019], we found that simply sampling policies and saving their cumulative associated rewards for each rollout  $r \in \{1, 2, \dots, K\}$ —then initializing with the medoid policy parameters corresponding to the highest reward—performed competitively against baselines. The full method is described in Algorithm 2. Following these considerations, we found that in few-shot settings ( $K = 10$  trajectory rollouts) CAML performs favorably to alternatives.

---

#### Algorithm 2 Training with CAML ( $K$ -medoids version)

---

**Require:** Distribution over entity types  $p(\mathcal{T})$ , initial policy parameters  $\theta$ , episode length  $T$   
**Require:** Number of medoids  $k$ , batch size to cluster with  $n$

- 1: **for** initial iterations  $1, \dots, n$  **do**
- 2:   Sample entity type  $\mathcal{T}_i \sim p(\mathcal{T})$  and initialize new policy  $\pi_\theta^i$
- 3:   Obtain and save updated  $\pi_{\theta'}^i$  and  $\tau_i'$  through Steps 3, 4, 5 in Algorithm 1
- 4: **end for**
- 5: Compute distance matrix  $D$  with pairwise distances on the saved  $\tau_i'$  using Algorithm 1.
- 6: Perform  $K$ -medoids clustering on  $D$ , saving set of  $k$  corresponding medoid policies  $\Pi_\theta^k$  and trajectories  $\tau_{1:N}^k$ . Save trajectories  $\tau^k \in \tau_{1:N}^k$  to compute next iteration of occupancy measures.
- 7: **for** iterations  $n + 1, n + 2, \dots$  **do**
- 8:   Sample entity type  $\mathcal{T}_i \sim p(\mathcal{T})$  and randomly sample medoid policy  $\pi_\theta \in \Pi_\theta^k$
- 9:   Obtain and save updated  $\pi_{\theta'}^i$  and  $\tau_i'$  through Steps 3, 4, 5 in Algorithm 1
- 10:   **if** number of saved policies equals  $n$  **then**
- 11:     Repeat steps 5 and 6, updating  $k$  medoid policies and trajectories.
- 12:   **end if**
- 13: **end for**

---

## 5 Experiments

In our experimental evaluation, we validate if: (1) clustering on our occupancy measure estimate reasonably characterizes divergence of policies over time, (2) our environments suggest that personalization is necessary for fast adaptation, and (3) building on these metrics can inspire effective few-shot learning algorithms. To focus on the effect of our algorithm we use vanilla policy gradient (VPG) for both meta-learning and evaluation.

### 5.1 Matching policy divergence to training environments

Towards evaluating both the dynamics of our testbed and the strength of our policy divergence distance metric, we first sought to see if clustering on trained VPG trajectories could recover their



Figure 2: **Policy divergence over time.** T-SNE visualization of policy divergences over 40 training updates. Colors denote original latent group. Computing pairwise distances using our divergence metric leads to noticeable grouping over time.

corresponding initial environmental dynamics. To do so, we first initiated six latent entity types by remapping controls completely (e.g. the yellow particle in Figure 1). Afterwards, for each latent type, we generated four variants by further introducing uniform variance for given entity and action while preserving the defining cardinal direction (red particle in Figure 1). In total we ended up with 24 entity types. To elucidate divergence across optimal policies, for each type we then trained an individual VPG policy from scratch over 40 updates performing a batch parameter update every 10 iterations, and calculating estimated occupancy-measure-based distances according to Algorithm 1.

As observed in Figure 2, based on estimated occupancy measure, trained policies show increasing signs of neighboring with members of their respective groups over training updates. One takeaway is that our occupancy measure-based divergence is effective at measuring policy divergence, where directed through their respective environments, we can thus evaluate how different two policies are based on their trajectories. Accordingly, in the few-shot RL setting, one valid strategy for producing an optimal policy without having to train from scratch would be to identify prior expert policies trained on similar environments. Taking this one step further, because our occupancy measures directly derive from  $\pi(a|s)$  which is itself derived from a policy network’s parameters, this provides another interpretation into the effectiveness of finding nearby parameters in policy space.

## 5.2 Evaluating fast adaptation to personalized environments

Given training on a set support types  $\mathcal{T}^S$ , our learning algorithms should quickly perform well on unseen query type set  $\mathcal{T}^Q$ . In-line with  $K$ -shot RL, for each support type  $\mathcal{T}_i^S \in \mathcal{T}^S$ , our policies are allowed  $K$  rollouts for adaptation during each training iteration. Towards evaluation, for each learning algorithm we train with 100 iterations using its policy-specific training algorithm. All underlying gradients were calculated with VPG, updating with batch size 10. For evaluation, due to the on-policy nature of policy gradient, we first allow  $K$  rollouts to collect new samples from unseen type  $\mathcal{T}_i^Q \in \mathcal{T}$ . Allocating these trajectories as an initialization period, for each  $\mathcal{T}_i^Q$  we then evaluate fast adaptation by comparing rewards across policies for up to five updates with fine-tuning via VPG.

We compare CAML to five other methods: (1) pretraining a policy by randomly sampling from all support environments during training, similar to joint training Yu et al. [2017], (2) pretraining a policy in the same query environment that we use for evaluation, (3) pretraining a policy in a different environment, (4) training with Reptile, a meta-learning algorithm comparable to MAML in performance but much simpler in computation that performs  $K > 1$  steps of stochastic gradient

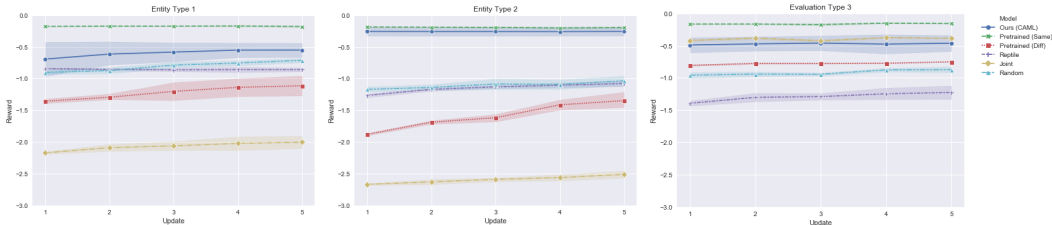


Figure 3: **Evaluation of learning algorithms on personalized environments.** CAML performs comparatively well across various evaluation entity types.

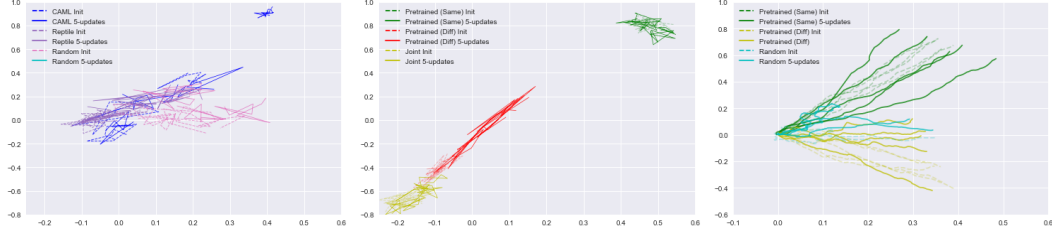


Figure 4: **Performance on 2D navigation.** Qualitative comparisons marking the ending location of particles at the end of episodes for the first evaluated type, with attention to performance after initialization and a few updates. Only the matched pretrained model and CAML reach positions noticeably closer to the target coordinate  $(1, 1)$  in 5 updates (left). Visualization of example trajectories taken by pretrained and random policies (right).

descent for randomly sampled support environments before directly moving initialization weights in the direction of the weights obtained during SGD Nichol et al. [2018], and (5) randomly initializing weights.

The results in Figure 3 suggest that CAML is able to quickly adapt to unseen entity types after one update during evaluation. We note that given three different environments, CAML consistently outperforms a support type-pretrained policy and a randomly-initialized policy, and also seems to outperform Reptile to varying degrees as well. The results also hint at differences in the environments dynamics. For example, entities with evaluation types 1 or 2 demonstrate instances where jointly trained and models pretrained on different environments initialize in poor states, suggesting that they require heavier personalization for success. On the other hand, evaluation type 3 represents a scenario where entity types may be representative of the larger population. While the matched pretrained policy serves as an upper bound in performance for all cases, we observe that the joint trained and unmatched pretrained policy model also seem to effectively transfer to the new environment.

Lastly we acknowledge the potential challenges of pretraining to personalized environments. Observing the ending coordinates of various policies in Figure 4, we note that in certain situations, such as when the query environment is not representative of a policy’s past experiences, methods such as joint training or pretraining with another environment may initiate detrimental behavior such as moving in the wrong direction.

## 6 Discussion and Future Work

We tackle meta-learning for fast adaptation in personalized environments, where environments may share the same reward, but differ dramatically in state transition probabilities. Using ideas from characterizing policy divergence, we introduce competitive improvements to a promising meta-learning algorithm paradigm, better initializing to a set of optimal policy parameters by better organizing past experiences into relevant clusters. While we believe further baselines and experiments are needed, our results also show preliminary support that meta-learning to minimize total distance between all potential optimal parameters may be sub-optimal when trying to adapt to diverse-enough environments. Our method builds on this original motivating curiosity, and suggests one possible alternative to better approach personalized reinforcement learning.



## References

- Andrew G Barto, Richard S Sutton, and Christopher JCH Watkins. Sequential decision problems and neural networks. In *Advances in neural information processing systems*, pages 686–693, 1990.
- Ignasi Clavera, Anusha Nagabandi, Ronald S. Fearing, Pieter Abbeel, Sergey Levine, and Chelsea Finn. Learning to adapt: Meta-learning for model-based control. *CoRR*, abs/1803.11347, 2018. URL <http://arxiv.org/abs/1803.11347>.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *CoRR*, abs/1703.03400, 2017. URL <http://arxiv.org/abs/1703.03400>.
- Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *CoRR*, abs/1606.03476, 2016. URL <http://arxiv.org/abs/1606.03476>.
- Kyle Hsu, Sergey Levine, and Chelsea Finn. Unsupervised learning via meta-learning. *CoRR*, abs/1810.02334, 2018. URL <http://arxiv.org/abs/1810.02334>.
- Leonard Kaufmann. Clustering by means of medoids. In *Proc. Statistical Data Analysis Based on the L1 Norm Conference, Neuchatel, 1987*, pages 405–416, 1987.
- Jianhua Lin. Divergence measures based on the shannon entropy. *IEEE Transactions on Information Theory*, 37(1):145–151, January 1991. ISSN 0018-9448. doi: 10.1109/18.61115. URL <http://ieeexplore.ieee.org/xpl/articleDetails.jsp?arnumber=61115>.
- Andreas Merentitis, Kashif Rasul, Roland Vollgraf, Abdul-Saboor Sheikh, and Urs Bergmann. A bandit framework for optimal selection of reinforcement learning agents. *CoRR*, abs/1902.03657, 2019. URL <http://arxiv.org/abs/1902.03657>.
- Anusha Nagabandi, Chelsea Finn, and Sergey Levine. Deep online learning via meta-learning: Continual adaptation for model-based RL. *CoRR*, abs/1812.07671, 2018. URL <http://arxiv.org/abs/1812.07671>.
- Andrew Y. Ng and Stuart J. Russell. Algorithms for inverse reinforcement learning. In *Proceedings of the Seventeenth International Conference on Machine Learning, ICML ’00*, pages 663–670, San Francisco, CA, USA, 2000. Morgan Kaufmann Publishers Inc. ISBN 1-55860-707-2. URL <http://dl.acm.org/citation.cfm?id=645529.657801>.
- Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *CoRR*, abs/1803.02999, 2018. URL <http://arxiv.org/abs/1803.02999>.
- Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B. Tenenbaum, Hugo Larochelle, and Richard S. Zemel. Meta-learning for semi-supervised few-shot classification. *CoRR*, abs/1803.00676, 2018. URL <http://arxiv.org/abs/1803.00676>.
- Susan M Shortreed, Eric Laber, Daniel J Lizotte, T Scott Stroup, Joelle Pineau, and Susan A Murphy. Informing sequential clinical decision-making through reinforcement learning: an empirical study. *Machine learning*, 84(1-2):109–136, 2011.
- Jake Snell, Kevin Swersky, and Richard S. Zemel. Prototypical networks for few-shot learning. *CoRR*, abs/1703.05175, 2017. URL <http://arxiv.org/abs/1703.05175>.
- Richard S Sutton, David A. McAllester, Satinder P. Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. A. Solla, T. K. Leen, and K. Müller, editors, *Advances in Neural Information Processing Systems 12*, pages 1057–1063. MIT Press, 2000. URL <http://papers.nips.cc/paper/1713-policy-gradient-methods-for-reinforcement-learning-with-function-approximation.pdf>.
- Wenhao Yu, Greg Turk, and C. Karen Liu. Multi-task learning with gradient guided policy specialization. *CoRR*, abs/1709.07979, 2017. URL <http://arxiv.org/abs/1709.07979>.