

## META REINFORCEMENT LEARNING

- Through previous experience (motivation of Meta RL)
- Data efficiency (with a few examples only)

- How might you get a machine to accomplish this task?

- Provide model of image formation, etc.
- Fewer human priors, more data-driven priors: Greater success.

## META-RL Problem Formulation and Examples

### META-LEARNING PROBLEM:

- In Supervised Learning:

Inputs:  $x$       Outputs:  $y$   
 $\vdash y = f(x; \theta) \uparrow$

DATA:  $\{ (x_i, y_i) \}_i^N$

- In Meta-supervised learning:

$D_{\text{train}}: \{ (x_i, y_i) \}_{i=1}^k$

Inputs:  $D_{\text{train}}$  and  $x_{\text{test}}$   
 $\uparrow$   
 $y_{\text{test}} = f(D_{\text{train}}, x_{\text{test}}; \theta)$

Outputs:  $y_{\text{test}}$   
 DATA:  $\{ D_i \}$   
 $\hookrightarrow$  diff. datasets (tasks)  
 $D_i: \{ (x_j, y_j) \}_j^N$

### EXAMPLE: FEW-SHOT CLASSIFICATION

- Given 1 example of 5 classes: we want to classify new examples. ( $x_{\text{test}}$ )

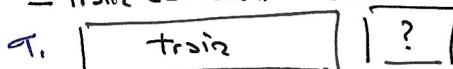
- Meta-Training:



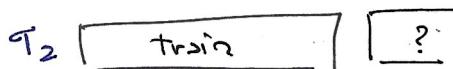
$\hookrightarrow$  which updates of function (the parameters)

#### META-TRAINING

- Train dataset #1: Cat-Bird -

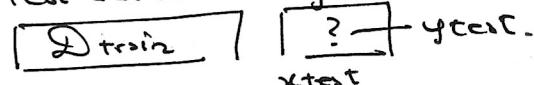


- Train dataset #2: flower-bike -



#### META-TESTING

- Test dataset: dog-ratter -



In one gradient step  $\delta$

DATA:  $\{ (s_t, a_t, r_t, s_{t+1}) \}$

- In Reinforcement Learning: Inputs:  $s_t$       Outputs:  $a_t$

$\vdash a_t = \pi(s_t; \theta) \uparrow$

- In Meta RL: Inputs:  $D_{\text{train}} \oplus s_t$       Outputs:  $a_t$

$k$ -rollouts  $\leftarrow$   
 from  $\pi$

$a_t = f(D_{\text{train}}, s_t; \theta) \uparrow$

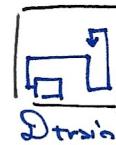
DATA:  $\{ D_i \}$

dataset of datasets collected for each task.

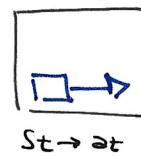
$\Rightarrow$  Design and Optimization of the  $f$  function  
 $\Rightarrow$  Collecting Appropriate Data (Learning to Explore - collect good Data)

### EXAMPLE: MAZE NAVIGATION

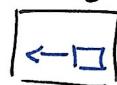
- Given a small amount of experience



$\Rightarrow$  Learn to solve the task



By learning how to solve many other tasks...



...

(other diff. mazes)

- For  $D_{train}$ :
  - (episodes)
    - $k$  rollouts from  $\pi \Rightarrow$  EPISODIC VARIANT
    - $t, \dots, k$  timesteps from  $\pi \Rightarrow$  ONLINE VARIANT.

### QUESTIONS:

#### 1. Meta-Learning vs Transfer Learning

- Optimize for good transfer performance (new task learnt differently)
  - $\Rightarrow$  Meta-Learning.
- Similar Motivations: both want to understand how to learn and transfer a representation that will help with learning new skills.
  - $\Rightarrow$  less formal for generalization process.
- TL: Accomplish this goal without any additional training.
- ML: Explicitly says that you can gather little more data on the new task to update the model (quickly)

### Method Classes

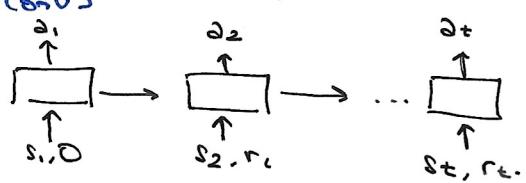
Design of the  $f$  function:

#### BLACK-BOX METHOD

$D_{train}, st \rightarrow at$

$\Rightarrow$  Recurrent Network  
(LSTM, NTH, Conv)

$$at = f(D_{train}, st; \theta)$$



$\Rightarrow$  diff compared to recurrent policy?

$\Rightarrow$  Hidden state maintained across episodes within a task.

- { +: general + expressive  
+: variety of design choices in Architecture  
-: complex model for complex task of learning  
-: Impractical data requirements.

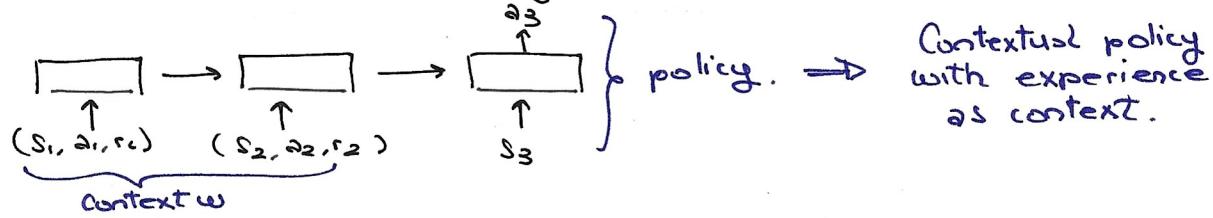
## PAPER: Simple Neural Attentive Meta-Learner

- Experiment: learning to visually navigate a maze
  - train on 1000 small mazes
  - test on held-out small mazes and large mazes.

## Digression: Connection to Contextual Policies.

- Contextual policy:  $\pi_\theta(a_t | s_t, w)$

↳  $w$ : stack location, walking direction, etc. (context)



What about goal-conditioned policies/value functions?

- Rewards: strict generalization of goals.

- Meta-RL: Adapt new tasks vs. Generalize new tasks

META-LEARNING  
(k-shot)

CLASSIC LEARNING  
(0-shot)

## OPTIMIZATION-BASED APPROACHES

Fine-Tuning  
[test-time]

$$\phi_j \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}_{\text{train}}(\theta)$$

↑  
pretrained  
parameters

↑  
Training data  
for new task.

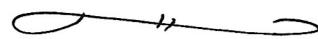
OUR METHOD

$$\min_{\theta} \sum_{\text{task } i} \mathcal{L}_{\text{test}}^i (\theta - \alpha \nabla_{\theta} \mathcal{L}_{\text{train}}^i(\theta))$$

↓  
in your  
meta-training  
set

↳ Embedded gradient  
descent in Meta-Learning

∴ Train over many tasks, to learn parameter vector  $\theta$  that transfers.

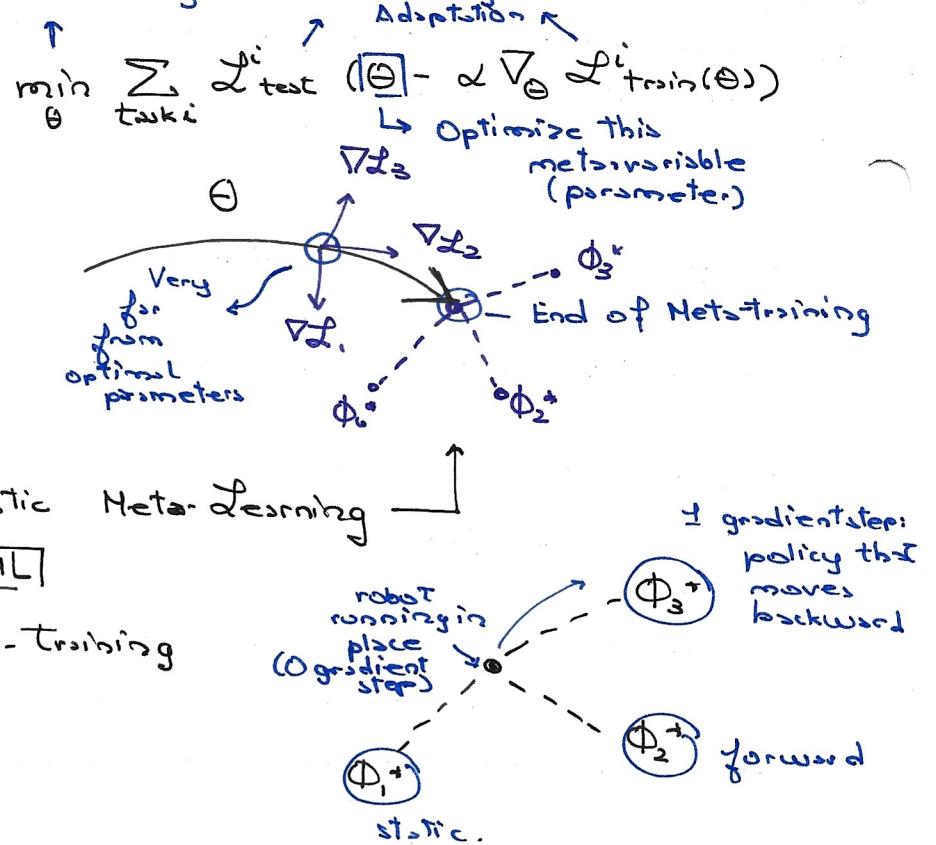


Suppose that:

- $\Theta$ : parameter vector being meta-learned
- $\Phi_i^*$ : optimal parameter vector for Task  $i$

{ ■ : Meta-Learning.  
--- : Learning / Adaptation }

Meta-Learning

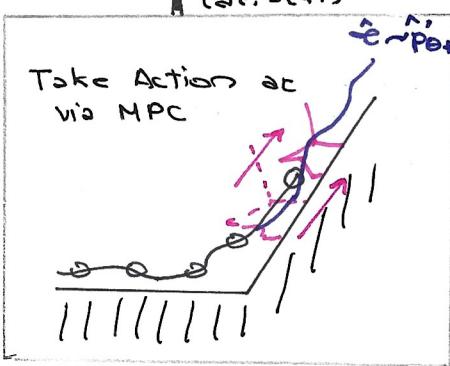
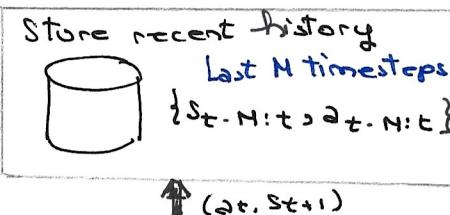


Example: At the end of meta-training

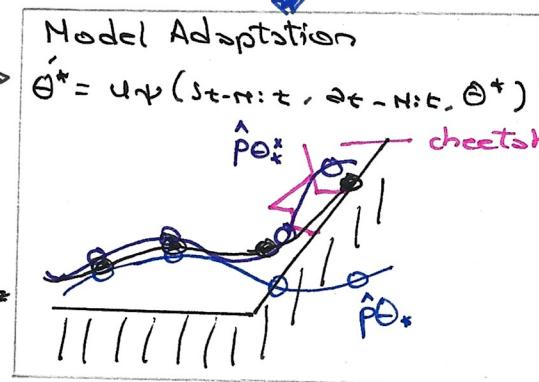
### MAML and Model-Based RL

- Adapt to diff. dynamics in environments.

### ONLINE ADAPTATION



Ex: On variable terrains  
META-TRAIN A PRIOR  $\Theta^*$



Ex: META-TRAIN on variable terrains

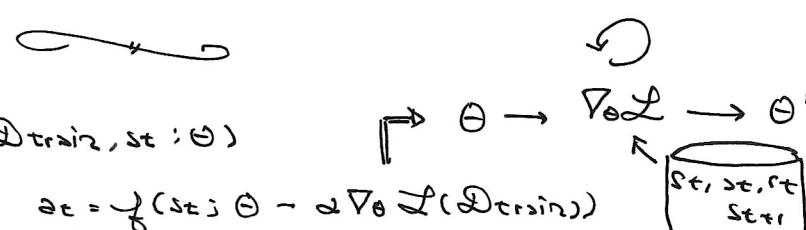
META-TEST with slope, missing leg, payload, calibration errors.

To summarize...

Black-Box Method :  $a_t = f(D_{\text{train}}, s_t; \Theta)$

Optimization-Based Method :  $a_t = f(s_t; \Theta - \alpha \nabla_{\Theta} \mathcal{L}(D_{\text{train}}))$

- +: inductive bias of SGD built-in
- +: model-agnostic (plug and play with existing archi.)
- : RL gradients & informative (policy gradients, Bellman errors)



## Learning to Explore

• How should we collect  $D_{train}$ ?

→ Coupling problem ↗ Execution  
Exploration ↗ Exploration

#0 Optimize for Exploration + Execution End-To-End w.r.t Reward

- (+) • simple
- leads to optimal strategy in principle

- (-) • Challenging optimization when exploration is hard

Example of hard Meta - RL Problem (Exploration):

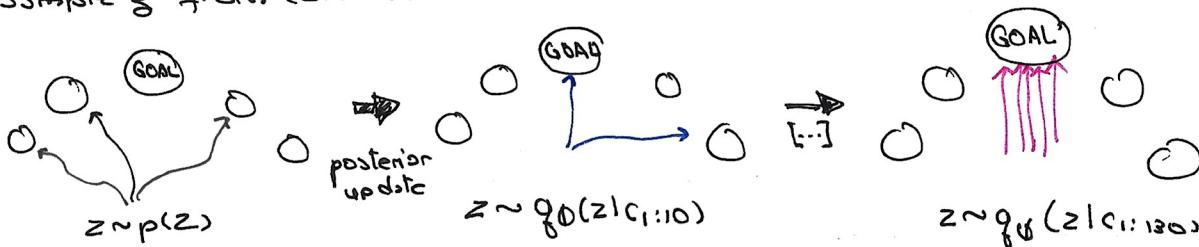
- { Learned cooking tasks in previous kitchens (META-TRAINING)  
 { Goal: quickly learn tasks in a new kitchen. (META-TESTING)
  - 1 exploration episode
  - 1 execution episode)

#1 Leverage Alternative Exploration Strategies.

a) Posterior sampling (Thompson sampling)

→ dots you've collected so far.

- Learn dist. over latent variable  $p(z)$ ,  $g(z|D_{train})$  and corresponding task policies  $\pi(a|s, z)$ .
- Sample  $z$  from current «posterior» and sample from policy  $\pi(a|s, z)$



⇒ Posterior sampling is bad when...

- (-) Goals far away and signs on wall that tells you the correct goal.

b) Use intrinsic rewards

c) Tasks dynamics and reward prediction

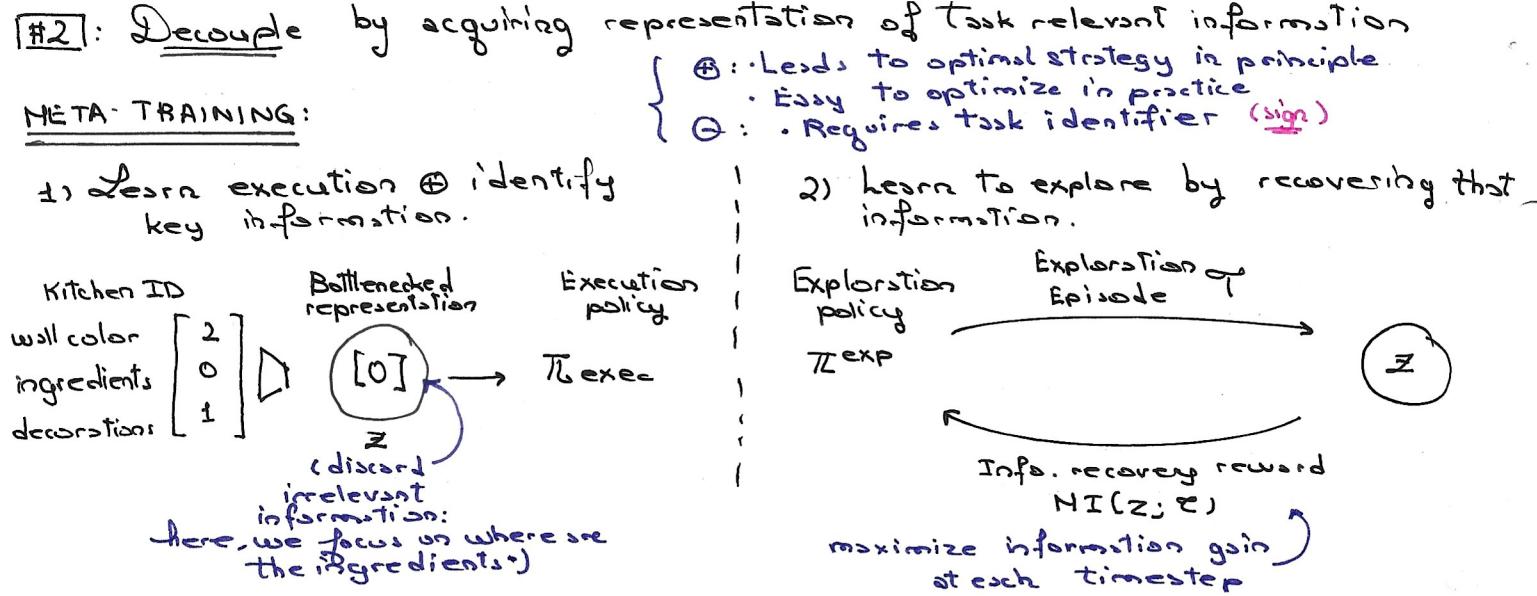
- Train model  $f(s, r | s, a, D_{train})$
- Collect  $D_{train}$  so that model is accurate.

Bad when...

- (-) Lots of distractors, or complex, high-dim state dynamics.

- (+) • Easy to optimize
- Many based on principled strategies

- (-) • Suboptimal by arbitrarily large amount in some environments.



META-TESTING:

$$\left\{ \pi^{\text{exp}} \xrightarrow{\text{Exploration Ep. } \mathcal{T}} \pi^{\text{exec}} \right\}$$

Decoupled Reward-free Exploration and Execution in Meta-RL (DREAM)  $\Rightarrow$  Consistent with End-to-End Approach.

Example: Sparse Reward 3D Visual Navigation Problem

- Task: go to the (key / block / ball) - color specified by the sign
- Agent starts on other side of barrier, must walk around to read the sign
- Pixel observations (80x60 RGB)
- Sparse Binary Reward

To summarize...

- End-to-End
- Alternative Strategies
- Decoupled Exp. and Exec (DREAM)

## Challenges and Latest Developments.

1. Adopting to entirely new tasks. (meta-world benchmark?)
2. Robustness to out of distribution tasks.
3. Where do the tasks even come from?
4. Better RL Algorithms.