

Replication study of MAML in Reinforcement Learning

Paper: Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks by Chelsea Finn et. al

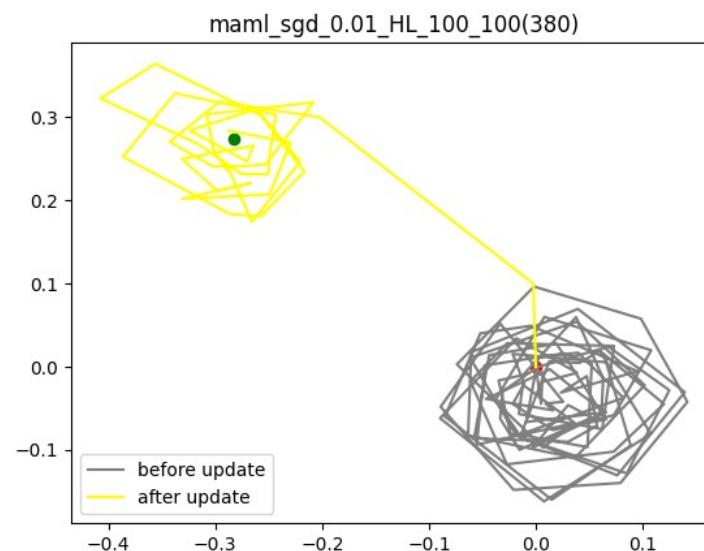
By André Henkel
Sommerterm 2020
University of Ulm
Institut für Neuroinformatik
Due Date: 18.08.2020

Scope of study

- Replication study of MAML in Reinforcement Learning
- Using 2D Navigation task
- Replicating the results of the paper
- Testing different setups
- Discuss found results and differences

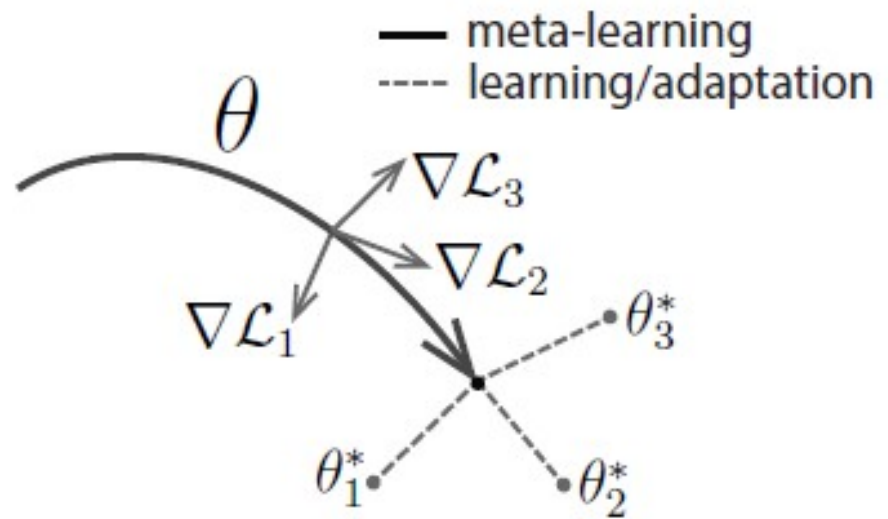
Why Meta-learning?

- Adapt to a new task quickly
- Few experience and updates
- For 2D Navigation:
 - Target perception through negative distance as reward
 - Observation includes the current position



MAML

- Model-agnostic meta-learning
- Train the model towards an initialization which adapts quickly to new tasks



MAML

- Initialization
- Inner-update
- Meta-update
- Loss of multiple tasks
- The paper asks three questions:
 - Quick adaptation
 - Multiple updates
 - Many domains

Algorithm 3 MAML for Reinforcement Learning

Require: $p(\mathcal{T})$: distribution over tasks

Require: α, β : step size hyperparameters

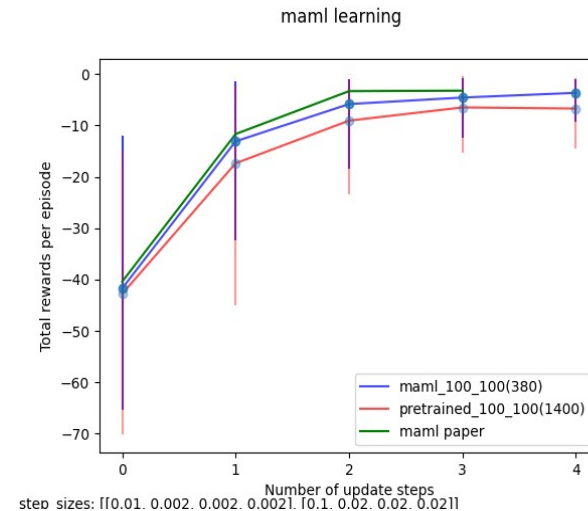
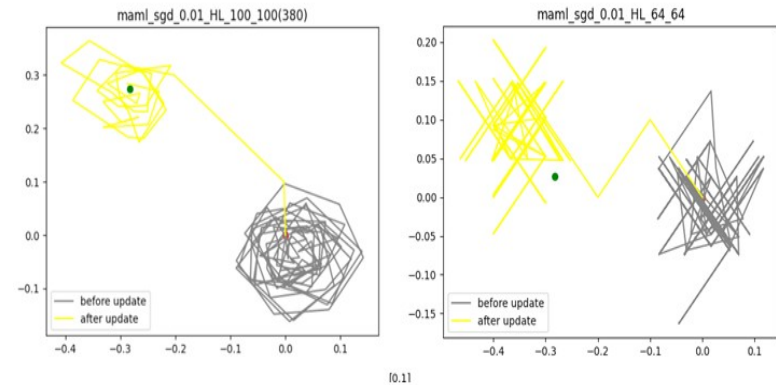
```
1: randomly initialize  $\theta$ 
2: while not done do
3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$ 
4:   for all  $\mathcal{T}_i$  do
5:     Sample  $K$  trajectories  $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H)\}$  using  $f_\theta$ 
      in  $\mathcal{T}_i$ 
6:     Evaluate  $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$  using  $\mathcal{D}$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation 4
7:     Compute adapted parameters with gradient descent:
       $\theta'_i = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ 
8:     Sample trajectories  $\mathcal{D}'_i = \{(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H)\}$  using  $f_{\theta'_i}$ 
      in  $\mathcal{T}_i$ 
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  $\mathcal{D}'_i$ 
      and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation 4
11: end while
```

Implementation

- Python, Pytorch
- Using TRPO for the meta-optimization
- Using most of the parameters as the authors
- Difference in step-sizes, policy variance and Neural Network output

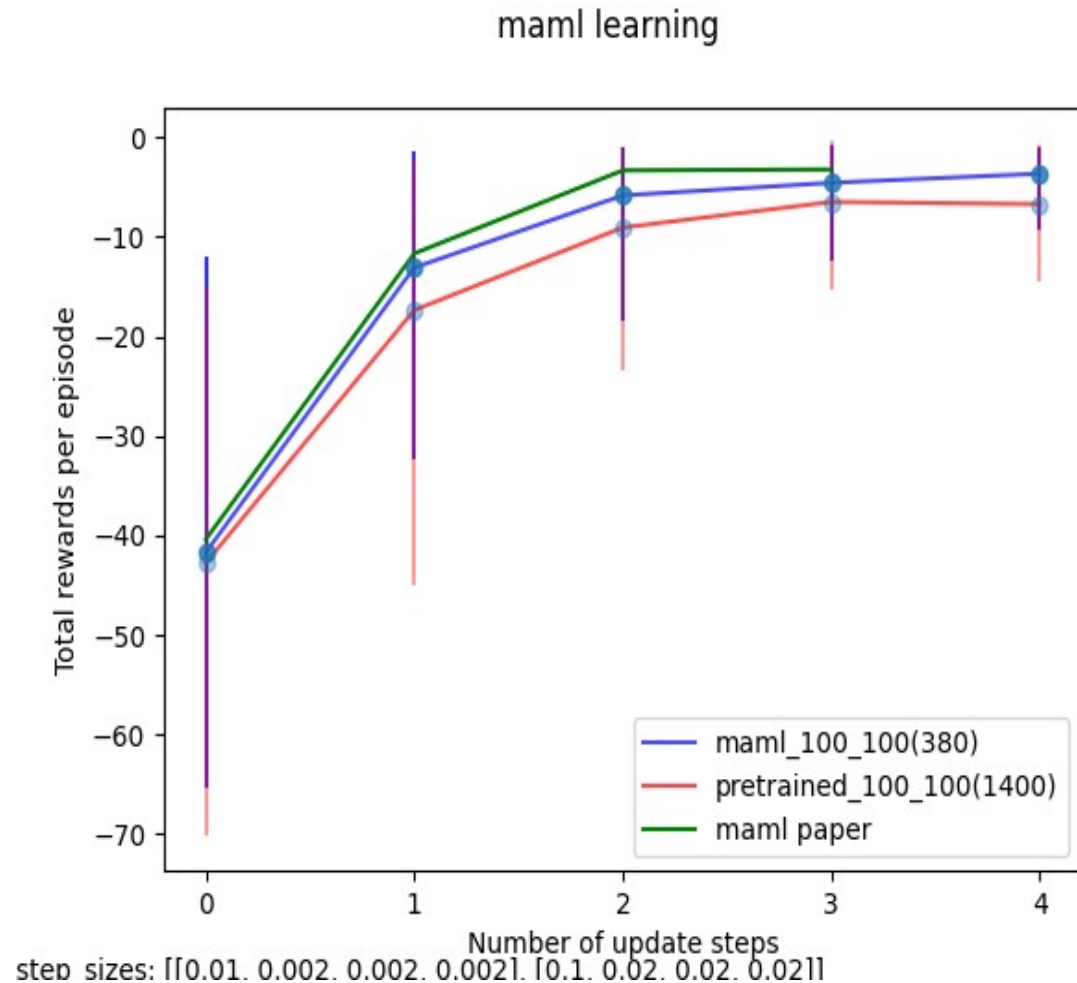
Testing and evaluation

- Display and test script
- Task distribution
- Setup for equal evaluation
- Tested setups:
 - Hidden layer sizes
 - Variance(exploration)
 - Learning rate
 - Meta-updates
 - Normal training
 - Out of meta range adaptation



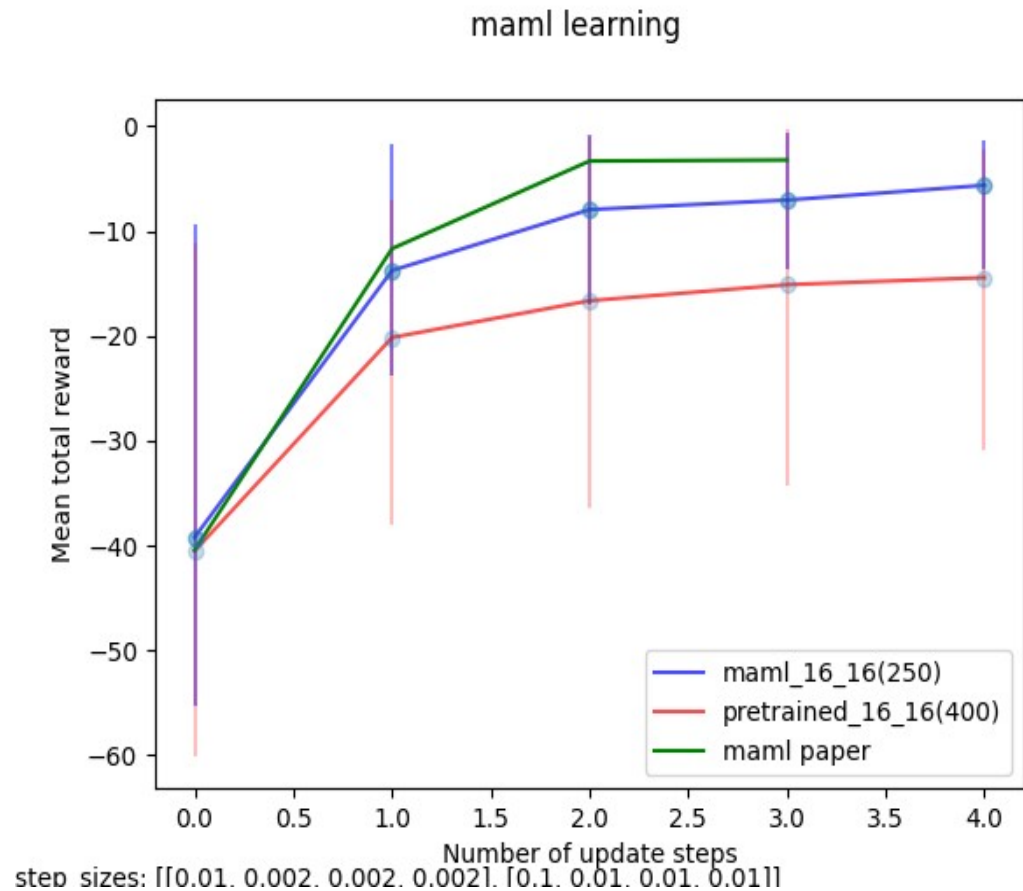
Adaptation comparison

- 100 units per HL
- 20 Tasks
- Errorbar(min/max)
- MAML paper comparison



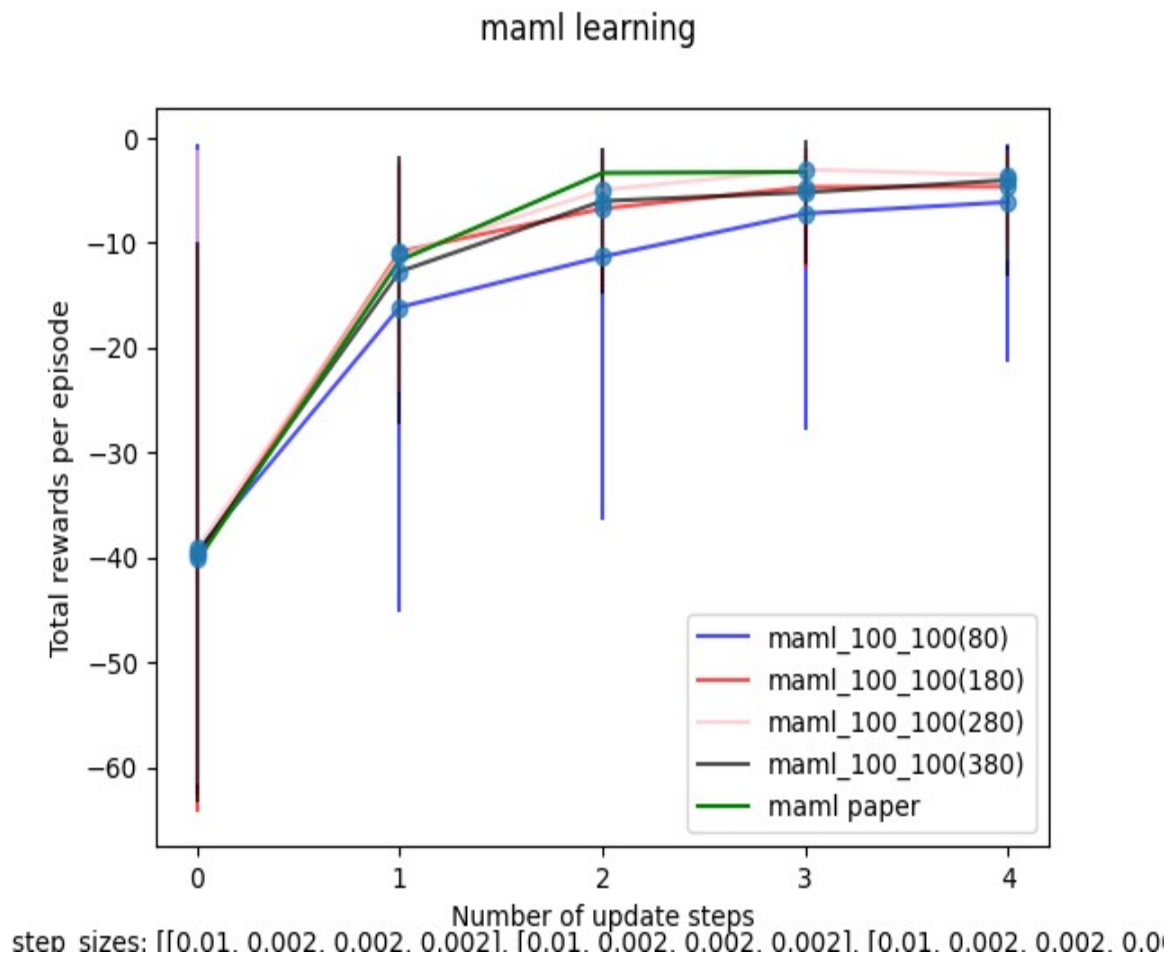
16 Units

- 16 Units per HL
- MAML still good
- Pretrained slightly worse



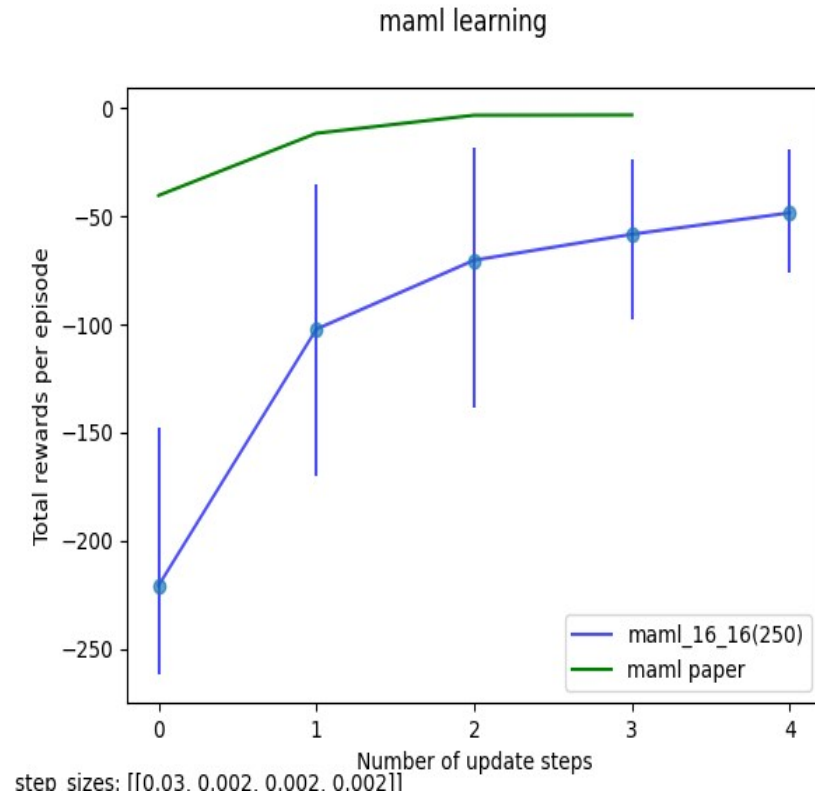
Effect of Meta-Updates

- 280 meta-updates in light orange performs best
- Same tasks comparison
- Just 80 updates has big min/max span and performs worse
- Empirically no more learning after ~380 updates



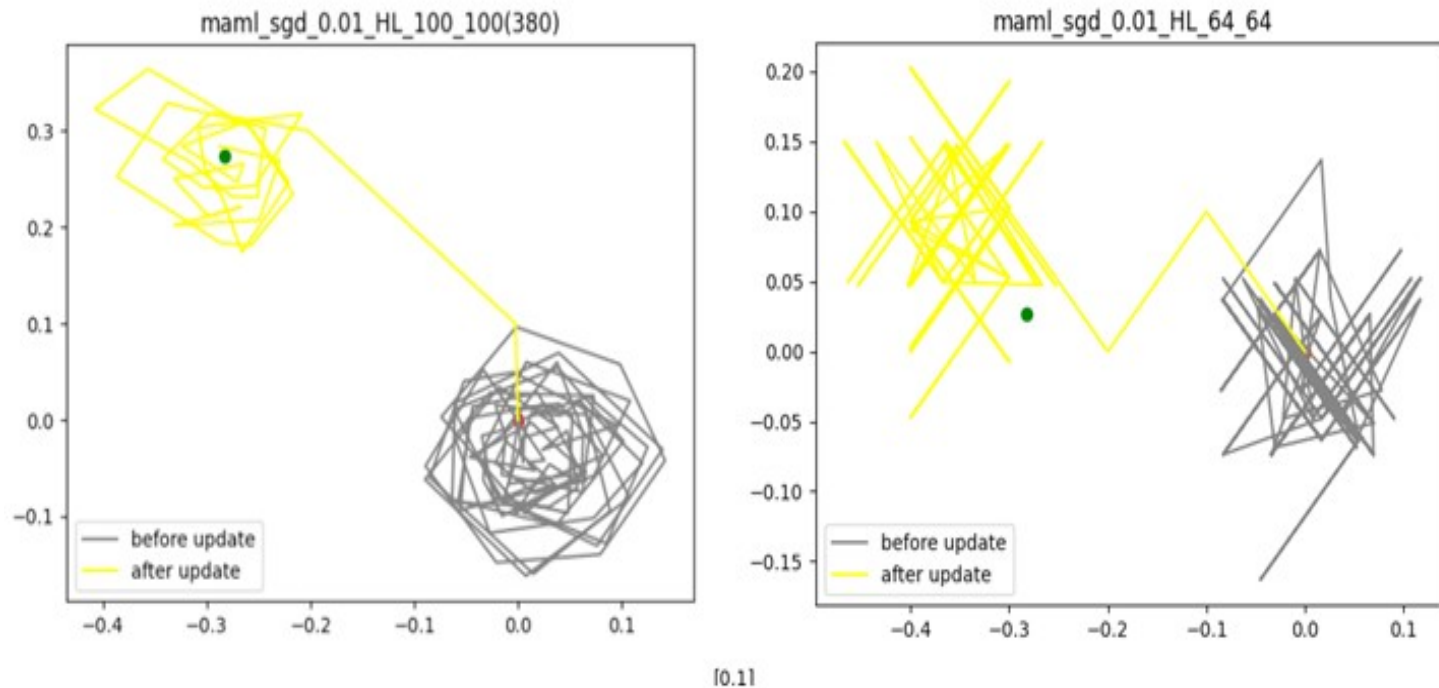
Out of range task performance

- Showing how a MAML pretrained model performs on tasks out of the meta-learned range
- Goal position $|1| - |2|$ each direction. Usually 0.5



Variance

- In official Github repo a variance with 1.0 and a FC output is taken
- In the replication study a sigma value of 0.1 and an tanh for the NN output is used.



Discussion

- Replicated results
- Initial paper questions could be answered
- Notable differences in implementation
- Other aspects and performances were tested