

# Meta-learnt priors slow down catastrophic forgetting in neural networks

Giacomo Spigler, CSAI, Tilburg University



## Abstract

Current training regimes for deep learning usually involve exposure to a single task / dataset at a time. Here we start from the observation that in this context the trained model is not given any knowledge of anything outside its (single-task) training distribution, and has thus no way to learn parameters (i.e., feature detectors or policies) that could be helpful to solve other tasks, and to limit future interference with the acquired knowledge, and thus catastrophic forgetting.

Here we show that catastrophic forgetting can be mitigated in a meta-learning context, by ex-posing a neural network to multiple tasks in a sequential manner during training. Finally, we present SeqFOMAML, a meta-learning algorithm that implements these principles, and we evaluate it on sequential learning problems composed by Omniglot and MiniImageNet classification tasks.

## Approach

We show that catastrophic forgetting is mitigated by extending a meta-learning framework to pre-train models on sequences of tasks, rather than individual tasks, to learn good, innate priors that are optimized for transfer and adaptability to tasks from a given distribution.

We propose to frame such type of sequential meta-learning as explicitly optimizing the performance of each task in a sequence both immediately after exposure to it, and at the end of the training sequence.

Specifically, we wish to optimize

$$\min_{\phi} \mathbb{E}_{\mathcal{T}_1, \dots, \mathcal{T}_n} \left[ \frac{1}{n} \sum_{i=1}^n L_{i, \text{test}}(\phi^n) + L_{i, \text{test}}(\phi^i) \right]$$

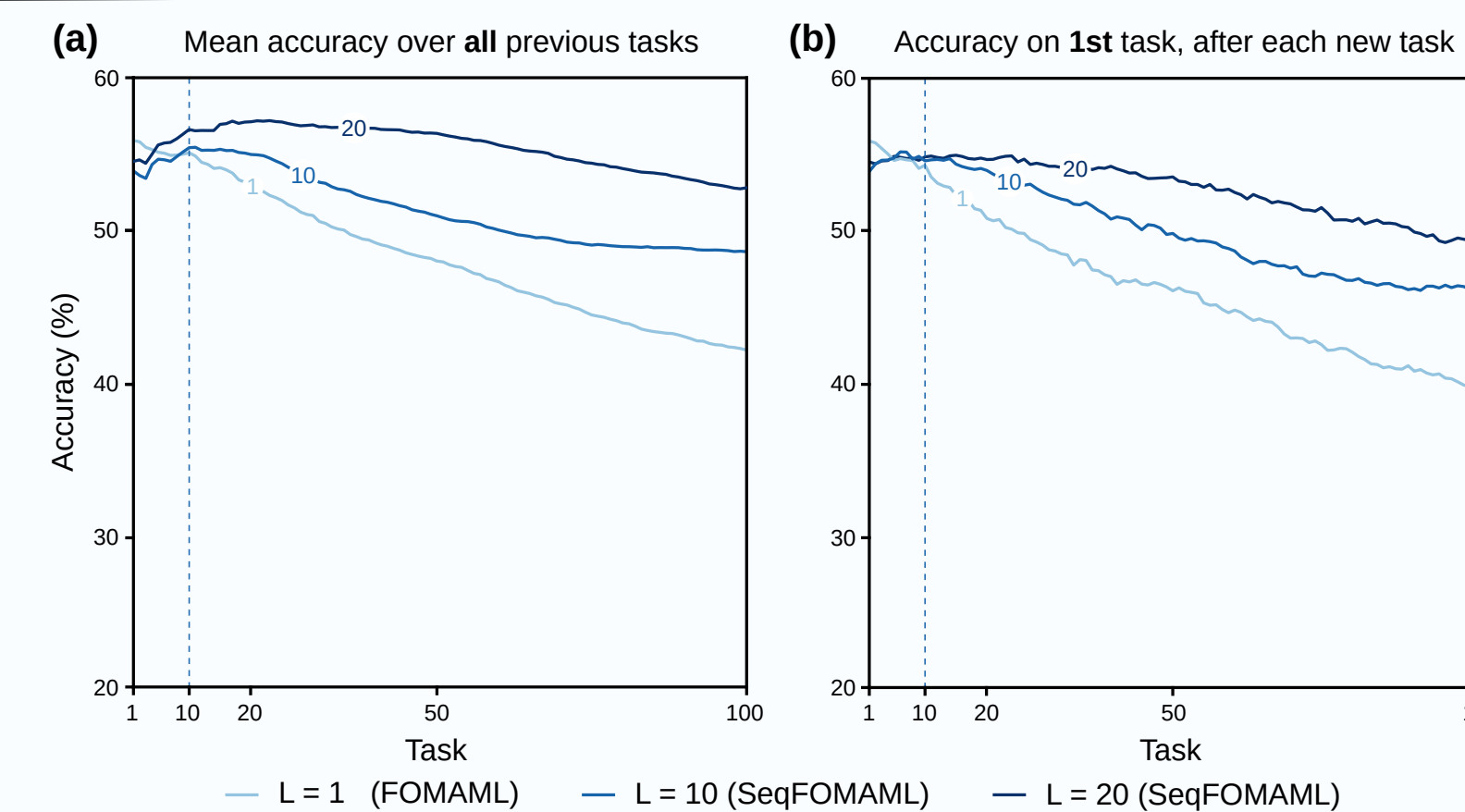
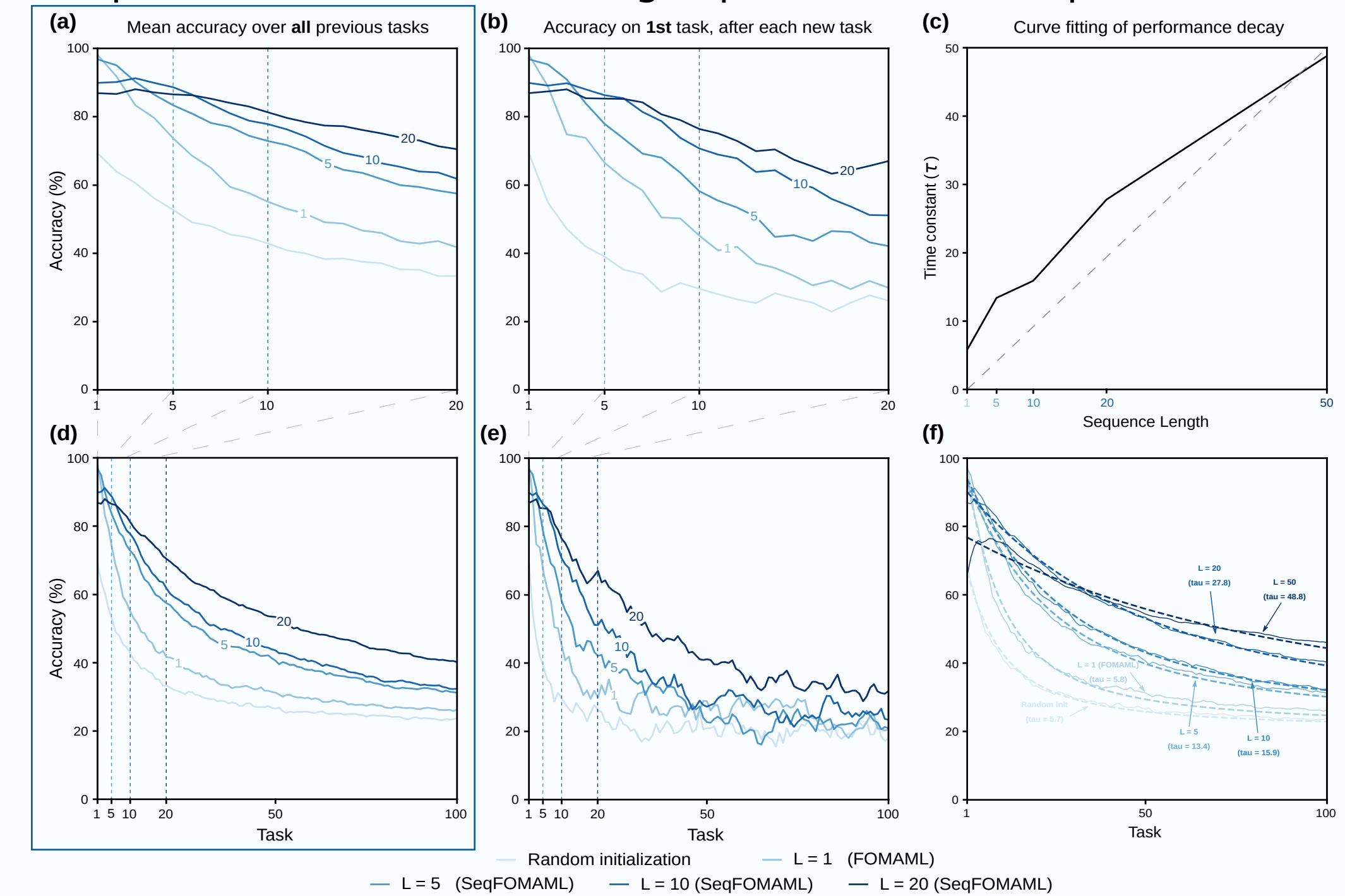
over a distribution of sequences of tasks of length  $n$ .  $\phi^i$  and  $\phi^n$  are the model parameters after training on task  $i$  and task  $n$  (end of training). The objective of the optimization problem is to find a set of initial model parameters such that the functional will yield a low loss when model parameters are fine-tuned by gradient descent on a novel sequence of tasks.

Gradients are then computed using (FO)MAML (Finn et al., 2017)

$$g_{\text{SeqMAML}} = \nabla_{\phi} \mathbb{E}_{\mathcal{T}_1, \dots, \mathcal{T}_n} \left[ \frac{1}{n} \sum_{i=1}^n L_{i, \text{test}}(\phi^n) + L_{i, \text{test}}(\phi^i) \right]$$
$$= \mathbb{E}_{\mathcal{T}_1, \dots, \mathcal{T}_n} \left[ \frac{1}{n} \sum_{i=1}^n L'_{i, \text{test}}(\phi^n) \frac{\partial \phi^n}{\partial \phi} + L'_{i, \text{test}}(\phi^i) \frac{\partial \phi^i}{\partial \phi} \right]$$
$$g_{\text{SeqFOMAML}} = \mathbb{E}_{\mathcal{T}_1, \dots, \mathcal{T}_n} \left[ \frac{1}{n} \sum_{i=1}^n L'_{i, \text{test}}(\phi^n) + L'_{i, \text{test}}(\phi^i) \right]$$

## Results

Incremental task learning, 5-way split-Omniglot (single-head network). Training is performed on sequences of 'L' tasks; testing is performed on sequences of 100 tasks.



5-way split-MiniImageNet  
(multi-head network).

The approach we presented has close connections with neuro-evolution approaches (see, e.g., works and reviews by Jeff Clune) and with the Baldwin effect (explored in the context of meta learning in Fernando et al., (2018)).

While the proposed system does not match the state-of-the-art performance of comparable methods (e.g., Javed & White (2019)), it shows that simply changing the training regime to account for the sequential presentation of tasks allows for a significant reduction in catastrophic forgetting.

**Note:** This page was made for screenshare-display at M2L. As such, it is less formal and lacks most details that would be in a traditional poster. Please refer to the arXiv paper (link below or QR code) for more details.

