# The effects of negative adaptation in Model-Agnostic Meta-Learning

Tristan Deleu, Yoshua Bengio

- The advantage of meta-learning is well-founded under the assumption that **the adaptation phase does improve the performance** of the model on the task of interest

- Optimization: maximize the performance after adaptation, **performance improvement is not explicitly enforced**

$$\min_{\theta} \mathbb{E}_{\mathcal{T} \sim p(\mathcal{T})}[\mathcal{L}(\theta'_{\mathcal{T}}; \mathcal{D}'_{\mathcal{T}})]$$

- We show empirically that performance **can decrease** after adaptation in MAML. We call this **negative adaptation**

- How to fix this issue? Ideas from **Safe Reinforcement Learning**