

SUMMARY

- In complex transfer learning scenarios, information contained only in the final parameters of a source model may be insufficient; a higher-level of abstraction is needed.
- We present a framework based on the idea that transfer learning can be achieved by leveraging information across similar learning processes, encoded in the geometry of the loss surface.
- We propose Leap, a lightweight meta-learner that scales beyond few-shot learning. Leap outperforms competing methods from meta-learning and transfer learning across a variety complex transfer learning scenarios.

SETUP: GRADIENT PATHS ON TASK MANIFOLDS

Given a learning objective f that consumes and input x and a target y and maps a parameterization θ to a scalar loss value, we define a learning process by the gradient update rule

$$\theta^{i+1} = \theta^i - \alpha \nabla f(\theta^i). \quad (1)$$

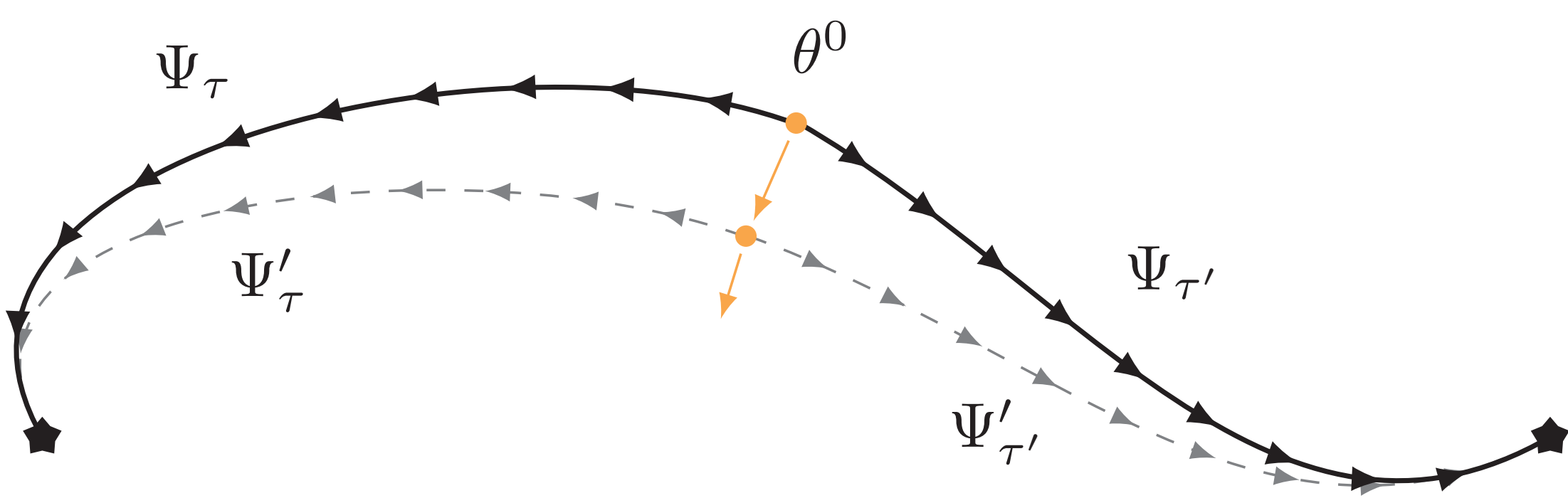
This process defines a *curve* γ on a task-specific manifold:

$$\text{Length}(\gamma) = \int_0^1 \sqrt{\langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle} dt. \quad (2)$$

- The choice of manifold determines how we view length. For simplicity, we define it as the loss surface: $\gamma(t) = (\theta(t), f(\theta(t)))$.
- A discrete approximation can be computed at negligible cost:

$$\text{Length}(\gamma) \approx \sum_{i=0}^{K-1} \|\gamma^{i+1} - \gamma^i\| = d(\theta^0). \quad (3)$$

Because the length of γ summarizes a learning process, we transfer knowledge by minimizing the *expected gradient path length* in Eq. 3 across a distribution $p(\tau)$ of tasks with unique paths Ψ_τ :



META LEARNING ACROSS TASK MANIFOLDS

- We focus on meta-learning a shared initialization θ^0 .
- Gradient paths cannot differentiate between good and bad local minima: we need a feasibility constraint.
- Given a second-best initialization ψ^0 , we aim to solve

$$\begin{aligned} \min_{\theta^0} \quad & F(\theta^0) = \mathbb{E}_{\tau \sim p(\tau)} [d_\tau(\theta^0)] \\ \text{s.t.} \quad & \theta^0 \in \Theta = \bigcap_{\tau} \{\theta^0 \mid f_\tau(\theta_\tau^*) \leq f_\tau(\psi_\tau^*)\}. \end{aligned} \quad (4)$$

LEAP

The feasibility constraint is costly to evaluate. Instead, we can use ψ^0 to generate baselines $\bar{\gamma}^i = (\psi^i, f(\psi^i))$ that guide θ^0 :

$$\bar{d}(\theta^0, \psi^0) = \sum_{i=0}^{K-1} \|\bar{\gamma}^{i+1} - \bar{\gamma}^i\|. \quad (5)$$

Leap uses $\bar{F}(\theta^0, \psi^0) = \mathbb{E}_{\tau \sim p(\tau)} [\bar{d}_\tau(\theta^0, \psi^0)]$ to obtain a sequence of incrementally demanding baselines that minimizes Eq. 4:

Theorem 1 (Pull-forward). Define a sequence of initializations $\{\psi_s^0\}_{s \in \mathbb{N}}$ by

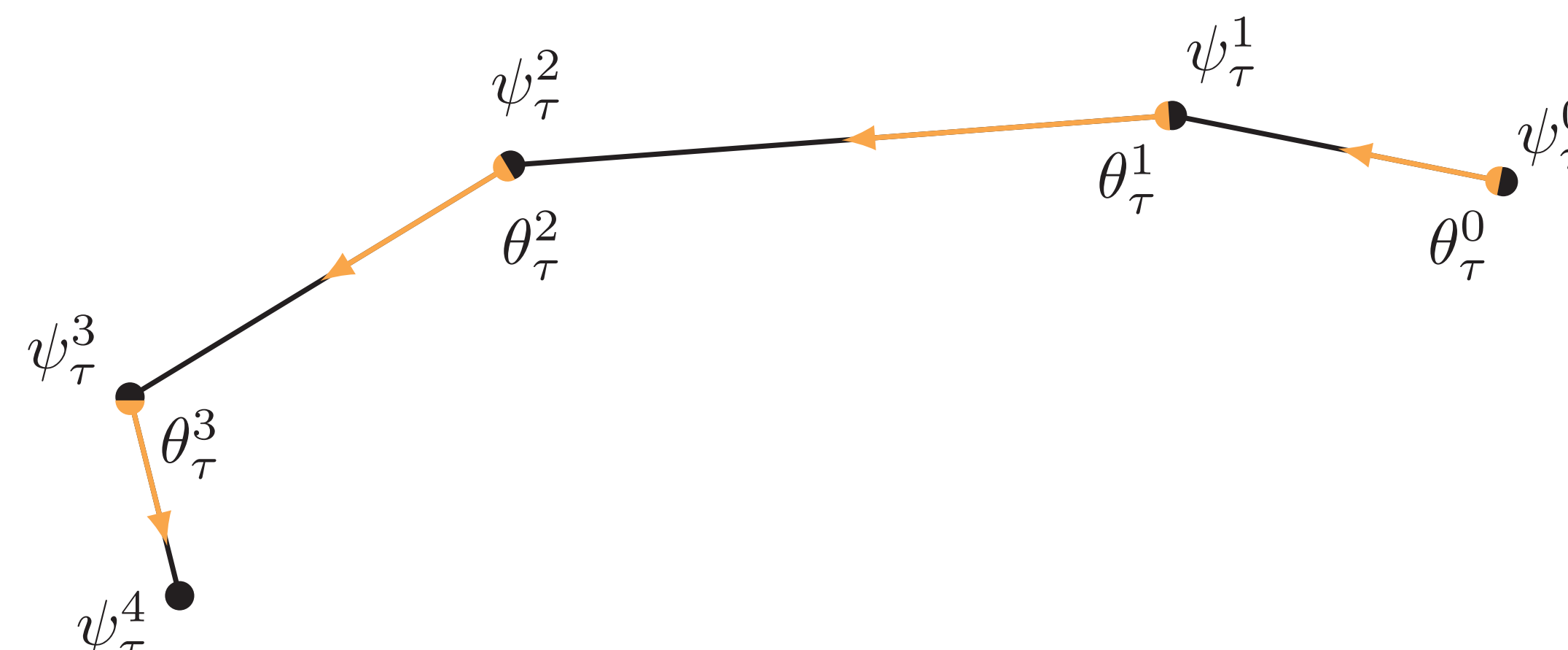
$$\psi_{s+1}^0 = \psi_s^0 - \beta_s \nabla \bar{F}(\psi_s^0, \psi_s^0), \quad \psi_0^0 \in \Theta. \quad (6)$$

For $\beta_s > 0$ sufficiently small, there exist learning rates schedules $\{\alpha_\tau^i\}_{i=1}^{K_\tau}$ for all tasks such that $\psi_{k \rightarrow \infty}^0$ is a limit point in Θ .

Crucially, the meta gradient can (approximately) be computed on the fly at negligible cost:

$$\nabla \bar{F}(\theta^0, \psi^0) \approx -\mathbb{E}_{\tau \sim p(\tau)} \left[\sum_{i=0}^{K_\tau-1} \frac{\Delta f_\tau^i \nabla f_\tau(\theta_\tau^i) + \Delta \theta_\tau^i}{(\|\bar{\gamma}_\tau^i - \gamma_\tau^i\|_2^p)^{2-p}} \right]. \quad (7)$$

Leap pulls the initialization forward along known gradient paths to find an initialization with minimal expected gradient path length that is guaranteed to perform as well as the baseline:



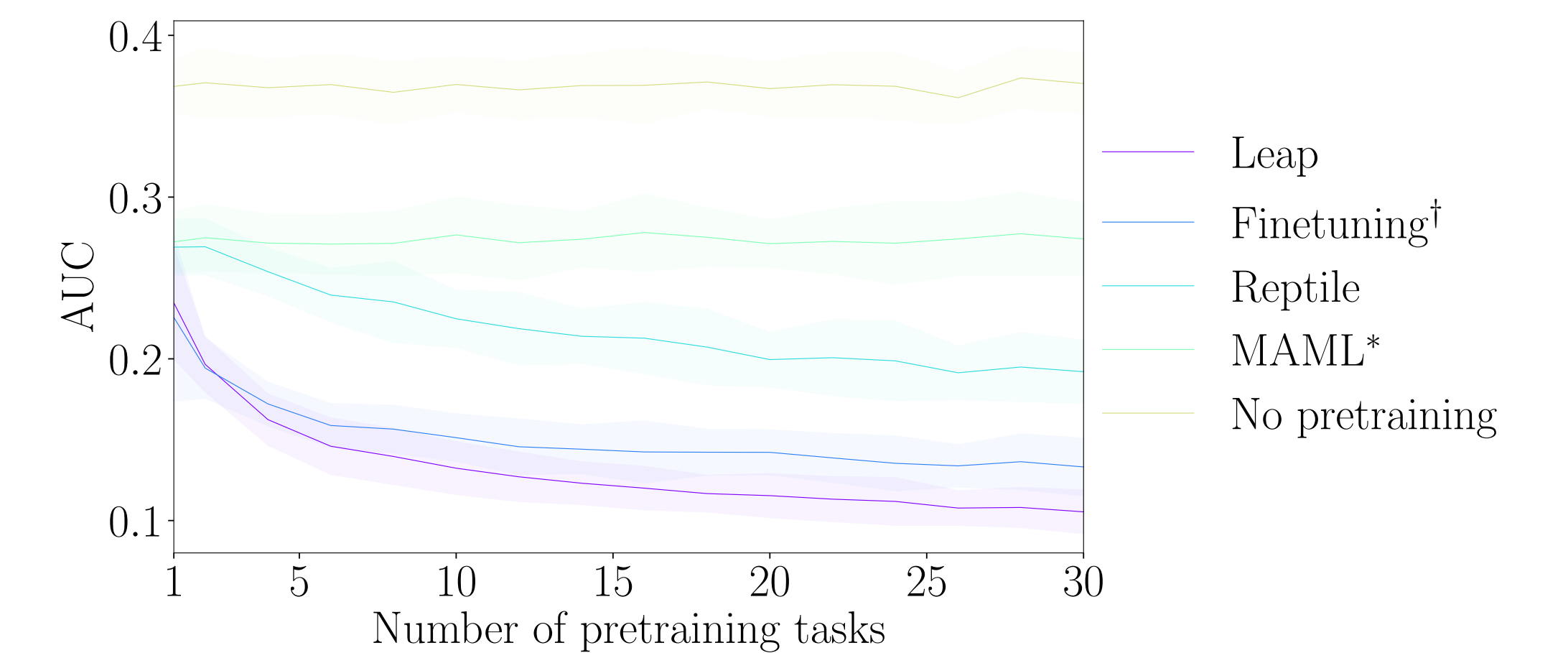
ALGORITHM 1: LEAP

Require: $\beta_s, p(\tau)$, $\tau = (f_\tau, u_\tau, p_\tau)$: distribution over tasks

- 1: randomly initialize θ^0
- 2: **while** not done **do**
- 3: $\nabla \bar{F} \leftarrow 0$: initialize meta gradient
- 4: sample task batch \mathcal{B} from $p(\tau)$
- 5: **for all** $\tau \in \mathcal{B}$ **do**
- 6: $\psi_\tau^0 \leftarrow \theta^0$: initialize task baseline
- 7: **for all** $i \in \{0, \dots, K_\tau - 1\}$ **do**
- 8: $\psi_\tau^{i+1} \leftarrow u_\tau(\psi_\tau^i)$: update baseline
- 9: $\theta_\tau^i \leftarrow \psi_\tau^i$: follow baseline (recall $\psi_\tau^0 = \theta^0$)
- 10: increment $\nabla \bar{F}$ using the pull-forward gradient (Eq. 7)
- 11: **end for**
- 12: **end for**
- 13: $\theta^0 \leftarrow \theta^0 - \frac{\beta}{|\mathcal{B}|} \nabla \bar{F}$: update initialization
- 14: **end while**

EXPERIMENTS

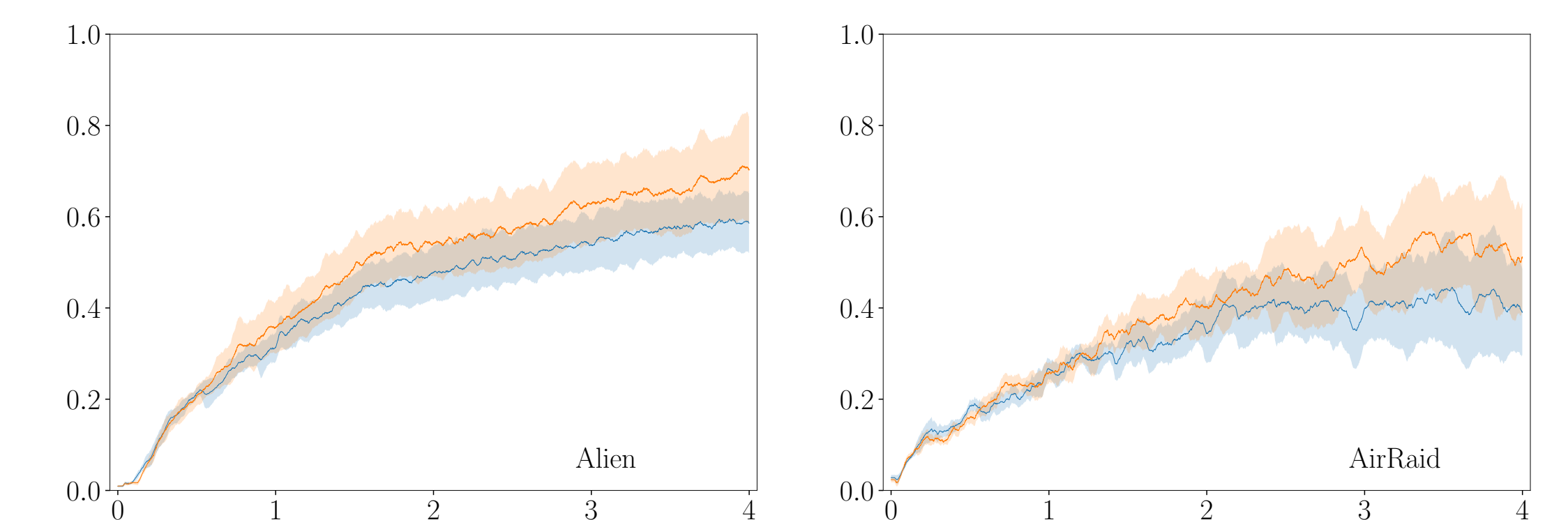
- Omniglot: each alphabet is a task, error AUC on test tasks:



- Multi-CV: a dataset is a task; mean normalized improvement:

	Leap	Finetuning	Progressive Nets [3]	HAT [4]
AUC	0.74	0.90	1.06	1.09
Test Error	0.89	1.20	0.97	1.15

- Atari: each game is a task; examples of Leap vs. random init:



References

- [1] C. Finn, P. Abbeel, S. Levine. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. ICML, 2017.
- [2] A. Nichol, J. Achiam, J. Schulman. On First-Order Meta-Learning Algorithms. arXiv:1803.02999, 2018.
- [3] A. Rusu, N. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, R. Hadsell. Progressive neural networks. J. Serrà, D. Suris, M. Miron, A. Kratzoglou. Overcoming catastrophic forgetting with hard attention to the task. ICML, 2018.
- [4] J. Serrà, D. Suris, M. Miron, A. Kratzoglou. Overcoming catastrophic forgetting with hard attention to the task. ICML, 2018.