

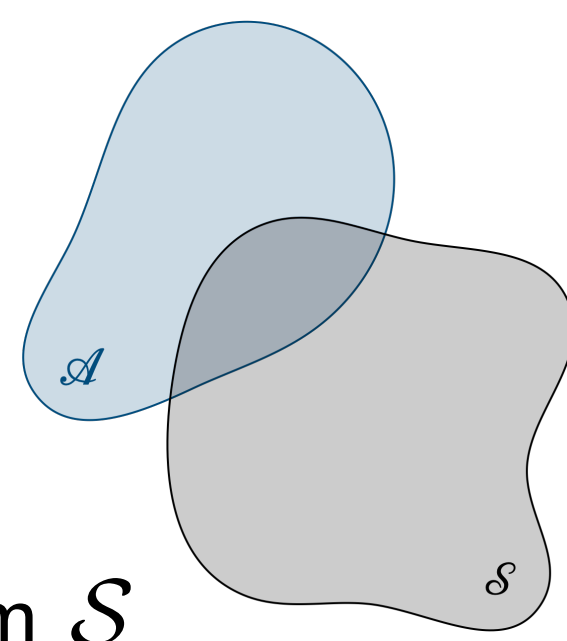
Background

- Active learning is a powerful paradigm for data selection
- Many fine-tuning tasks have lots of structure that can be identified by large pre-trained models and used to effectively select data for fine-tuning
- So far: the growing literature on fine-tuning has primarily focused on architectures and optimization algorithms

How can we select the *right* data for fine-tuning a large NN to a specific task?

Transductive Active Learning

- Sample space $\mathcal{S} \subseteq \mathcal{X}$ (train set)
- Target space $\mathcal{A} \subseteq \mathcal{X}$ (test set)
- Unknown function f over \mathcal{X}



Goal: Learn f within \mathcal{A} by sampling from \mathcal{S}

Can we exploit the (unknown) **latent structure** for fine-tuning?

- Embeddings $\phi(\cdot)$ generated by the NN capture the latent structure of f (e.g., neural tangent embs.)
- Approximate NN by a Gaussian process with kernel $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$, which encodes “similarity”
- The marginal variance $\sigma_n^2(\mathbf{x})$ of the Gaussian process after n samples is a proxy for the approximation error at \mathbf{x} after fine-tuning on these n samples

New goal: Reduce uncertainty $\sigma_n^2(\mathbf{x})$ at $\mathbf{x} \in \mathcal{A}$

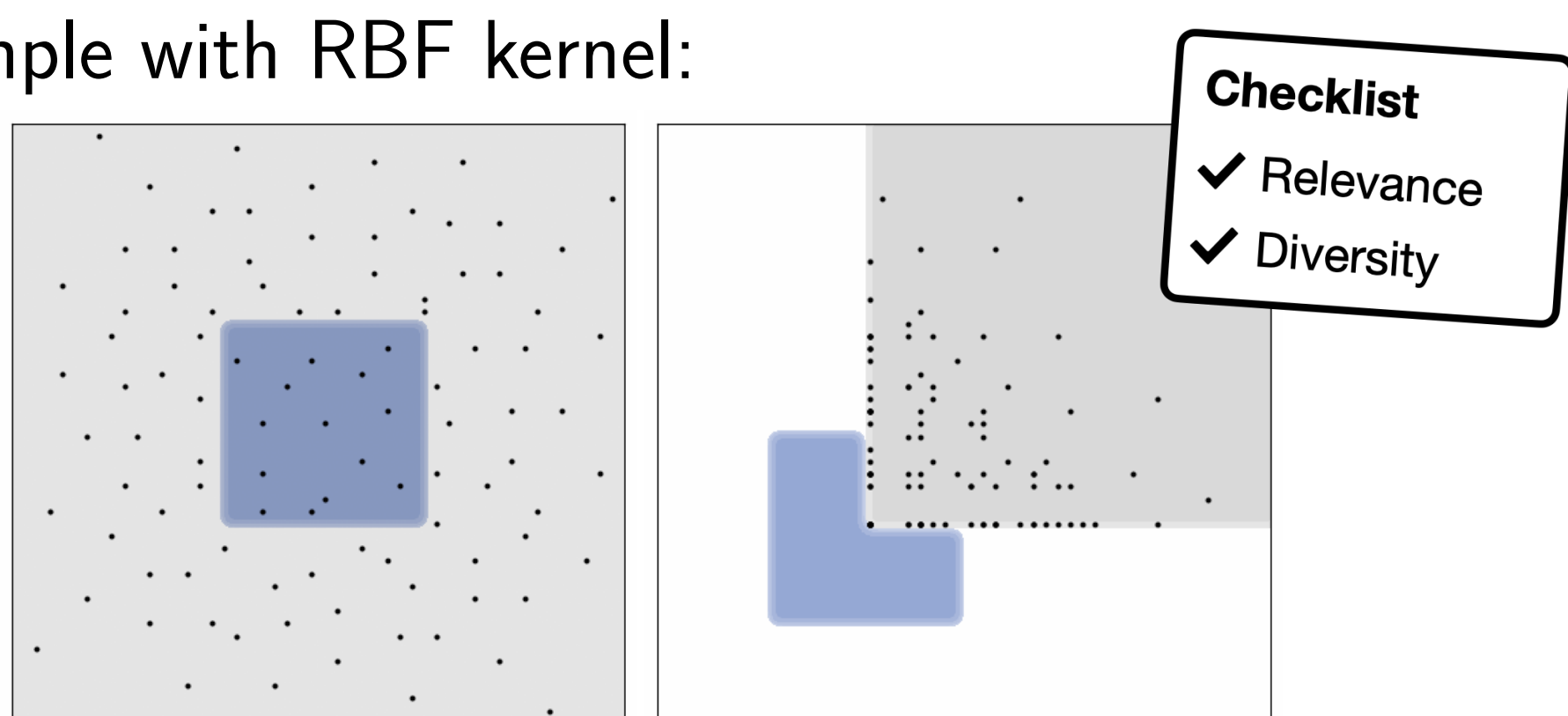
Information-based Transductive (Active) Learning

Proposal: Select samples to minimize posterior “uncertainty” within \mathcal{A}

Quantifying “uncertainty” by entropy yields **ITL**:

$$\begin{aligned} \mathbf{x}_n &= \arg \min_{\mathbf{x} \in \mathcal{S}} H[\mathbf{f}(\mathcal{A}) \mid D_{n-1}, (\mathbf{x}, y)] \\ &= \arg \max_{\mathbf{x} \in \mathcal{S}} I(\mathbf{f}(\mathcal{A}); (\mathbf{x}, y) \mid D_{n-1}) \end{aligned}$$

Example with RBF kernel:



Theory: Convergence Guarantees

How much can be learned about \mathcal{A} from \mathcal{S} ?

Generalization bound for ITL. For every $\mathbf{x} \in \mathcal{A}$:

$$\sigma_n^2(\mathbf{x}) \leq \underbrace{\text{Var}[f(\mathbf{x}) \mid \mathbf{f}(\mathcal{S})]}_{\text{irreducible}} + \underbrace{C \log(n)/\sqrt{n}}_{\text{reducible}}$$

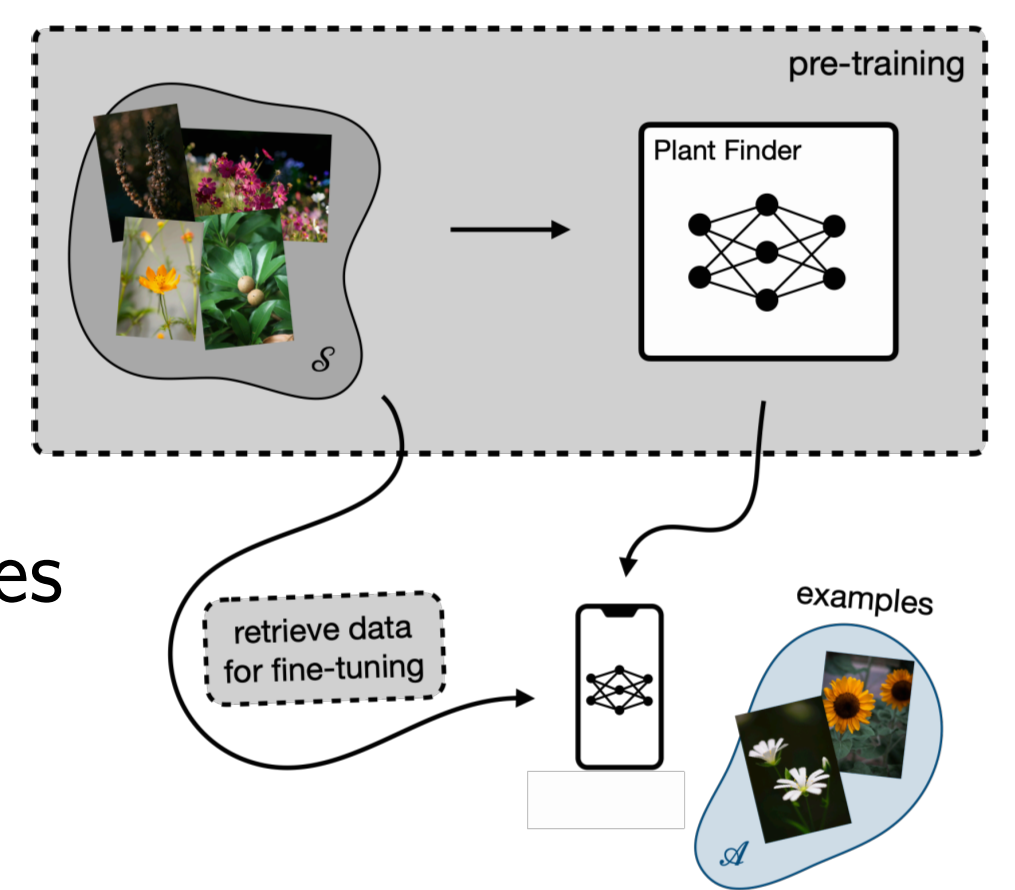
Approximation error bound for ITL. If $f \in \mathcal{H}_k(\mathcal{X})$ then for every $\mathbf{x} \in \mathcal{A}$ with probability $1 - \delta$:

$$|f(\mathbf{x}) - \mu_n(\mathbf{x})| \leq \beta_n(\delta) \left[\text{irreducible} + C \log(n)/\sqrt{n} \right]$$

where $\mu_n(\mathbf{x})$ is the prediction and $\beta_n(\delta)$ the CI width

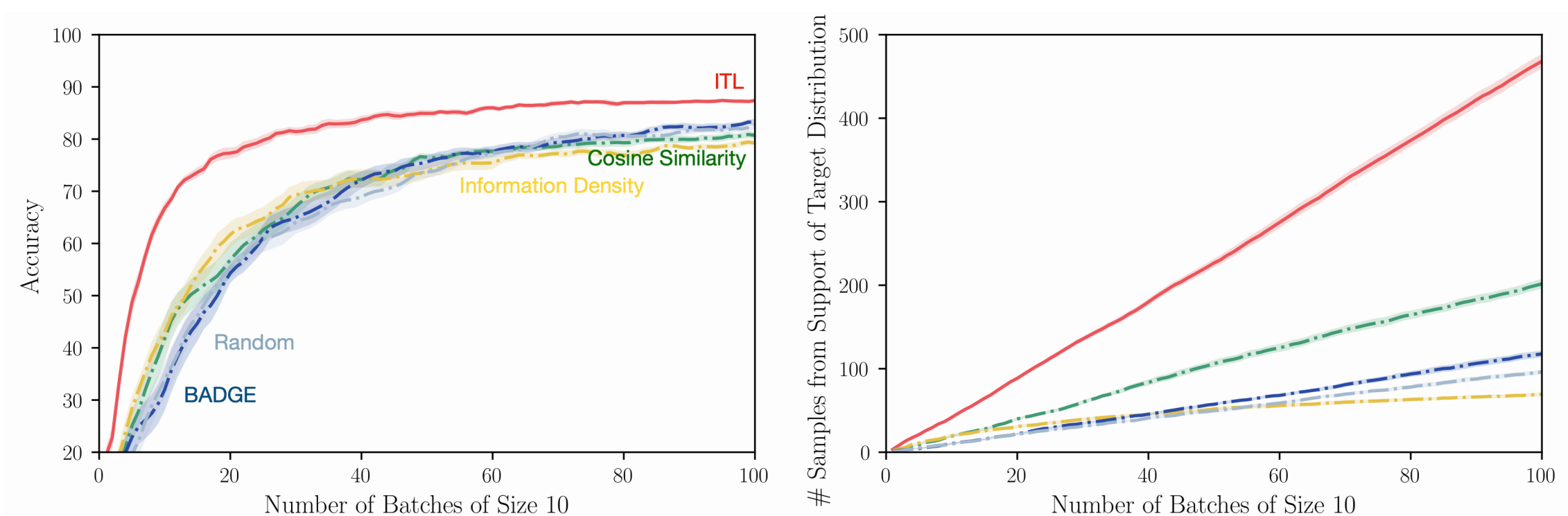
Practice

A motivating example: fine-tuning a plant classifier on a user’s local biome



Find informative (that is, relevant & diverse) examples for the user’s biome

Simplified example: fine-tuning on CIFAR-100



Cosine Similarity \triangleleft :
only relevance $\arg \max_{\mathbf{x} \in \mathcal{S}} \sum_{\mathbf{x}' \in \mathcal{A}} \mathbb{1}_{\cos(\phi(\mathbf{x}), \phi(\mathbf{x}')) > \tau} \mathbb{1}_{D_{n-1}}$

Information Density:
only relevance

BADGE:
only diversity

ITL:
relevance + diversity

ITL generalizes \triangleleft to query & batch sizes larger than 1

ITL synthesizes approaches to retrieval & coverage

Key Takeaways

- Retrieving the **right** examples for fine-tuning can lead to substantial performance gains
- Transductive active learning is a powerful paradigm for learning under resource constraints such as limited compute time & limited access
- ITL can be used as a simple drop-in replacement for random data selection