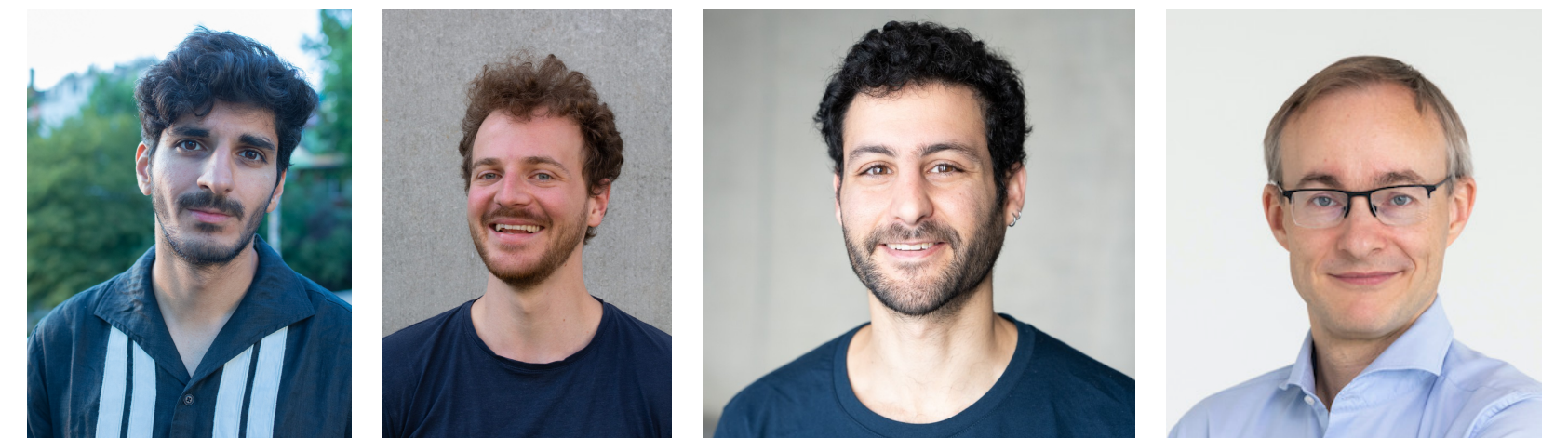


Active Fine-Tuning of Large Neural Networks

Jonas Hübötter

with Bhavya Sukhija, Lenart Treven, Yarden As, Andreas Krause



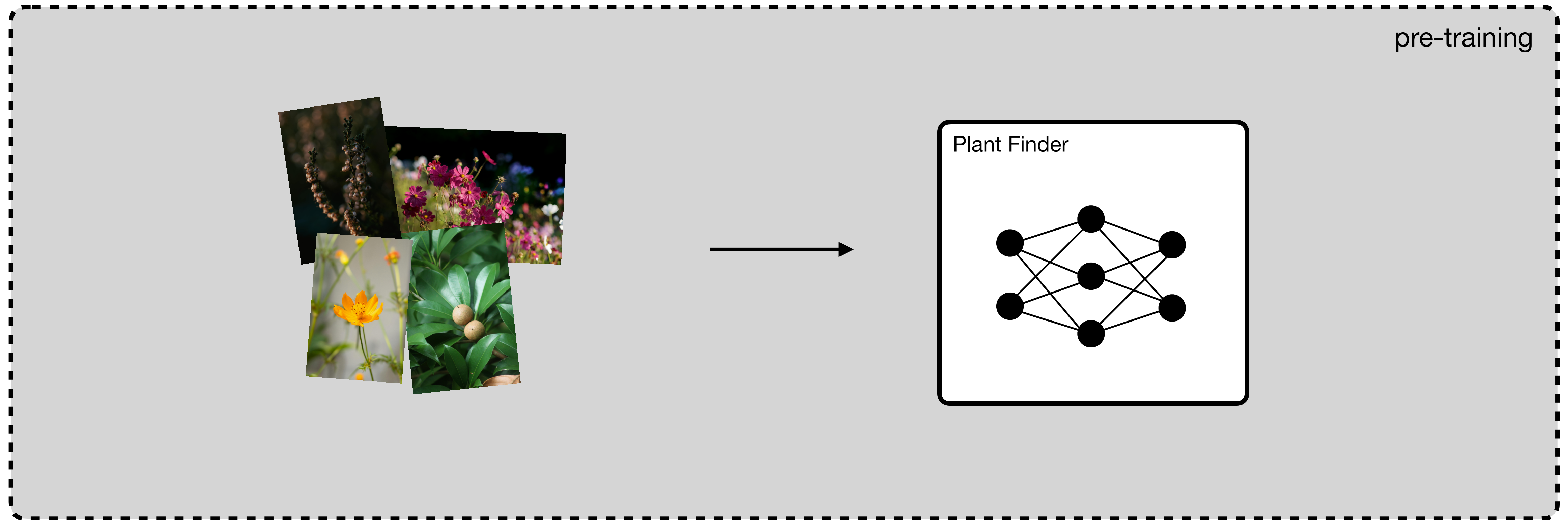
What I will *not* talk about: How to use a small dataset for fine-tuning

What I will *not* talk about: How to use a small dataset for fine-tuning

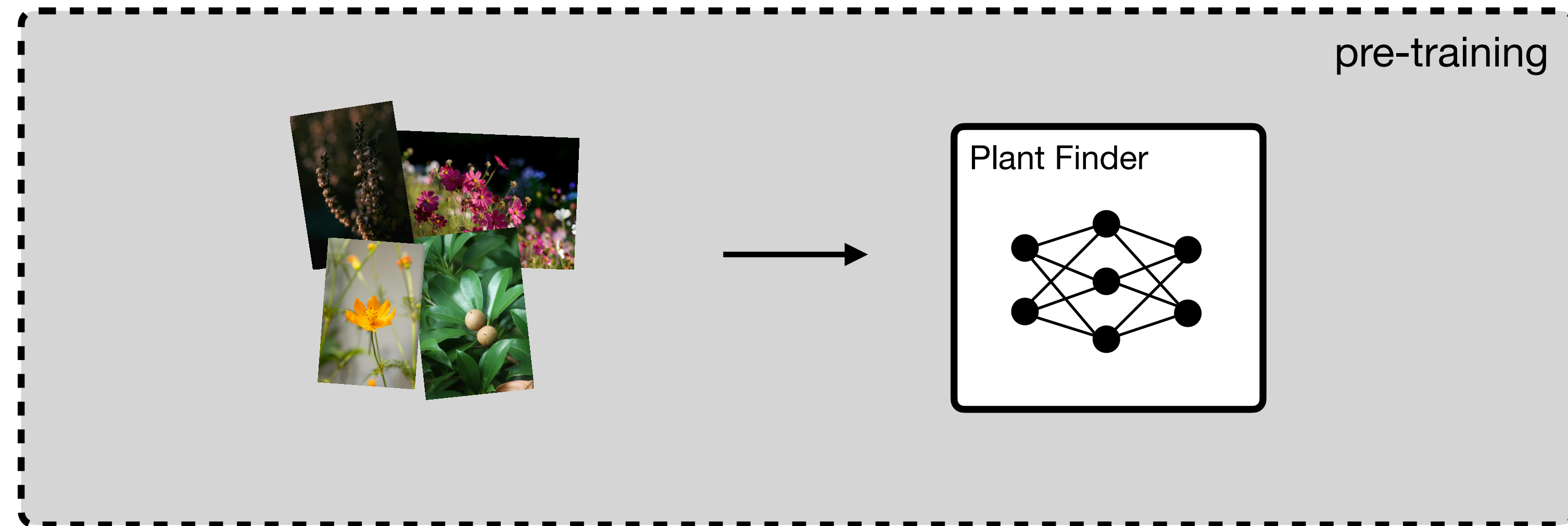
What I aim to convince you of:

Retrieving the **right** examples for fine-tuning can lead to substantial performance gains

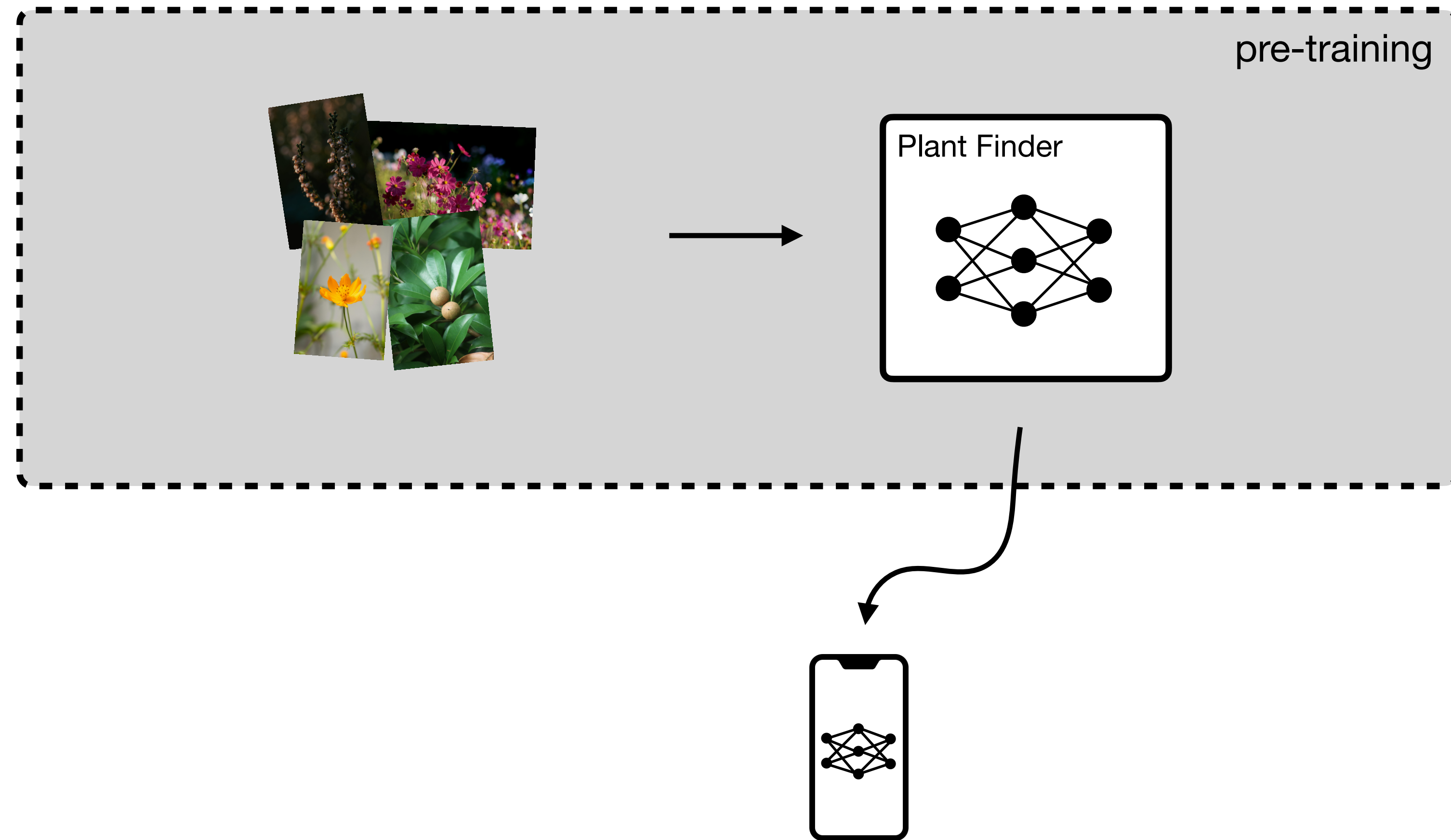
A Motivating Example



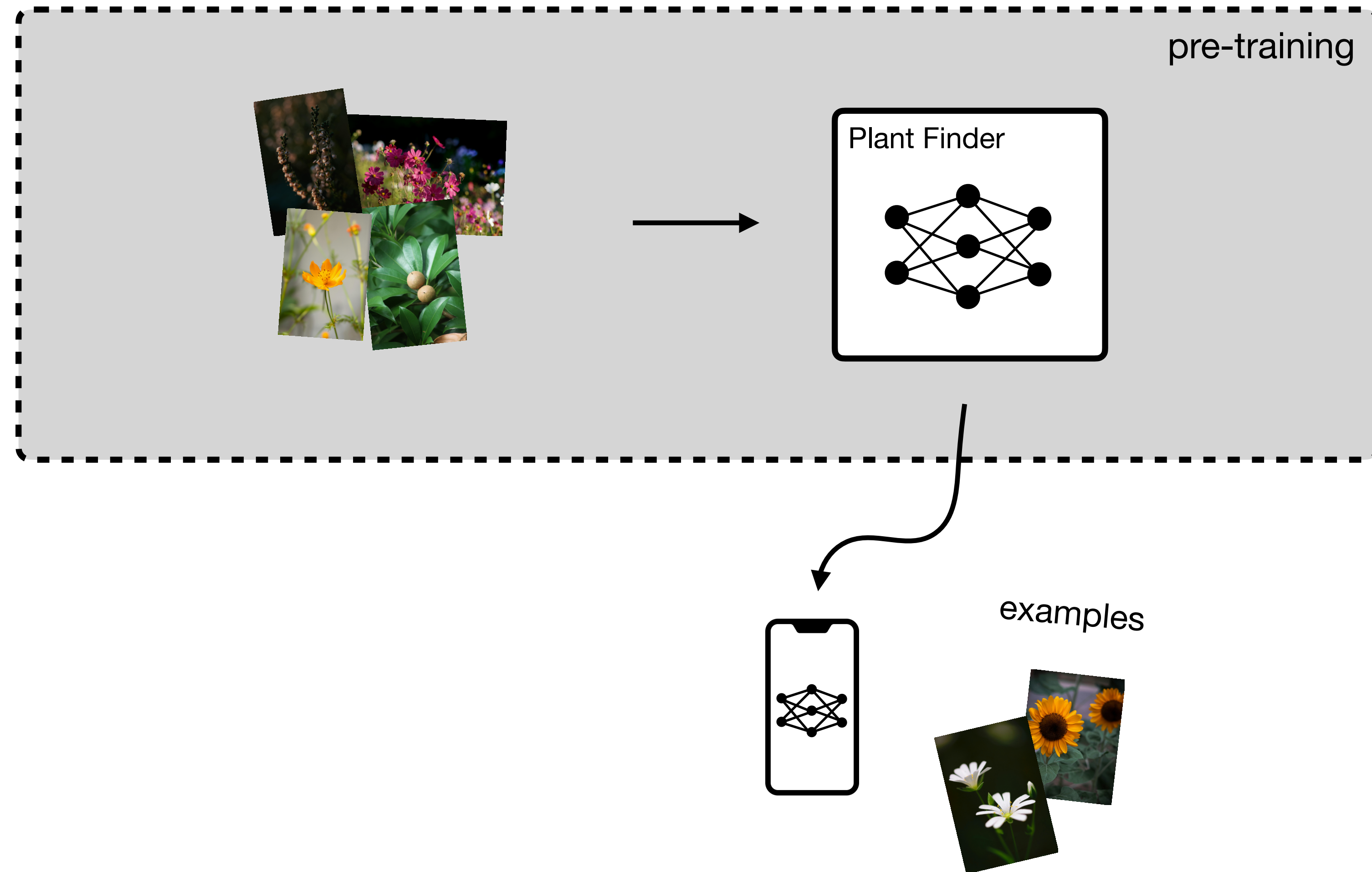
A Motivating Example



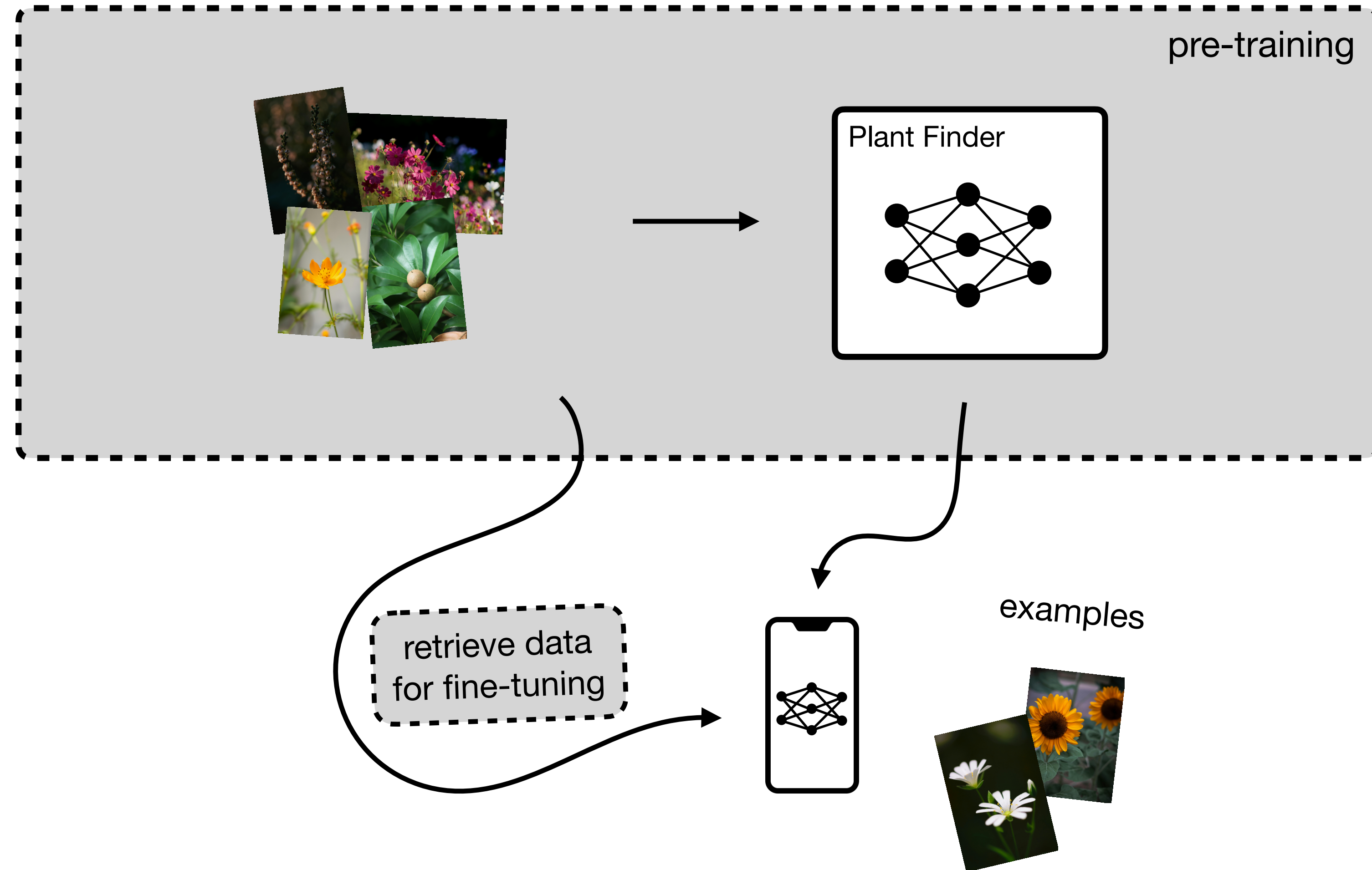
A Motivating Example



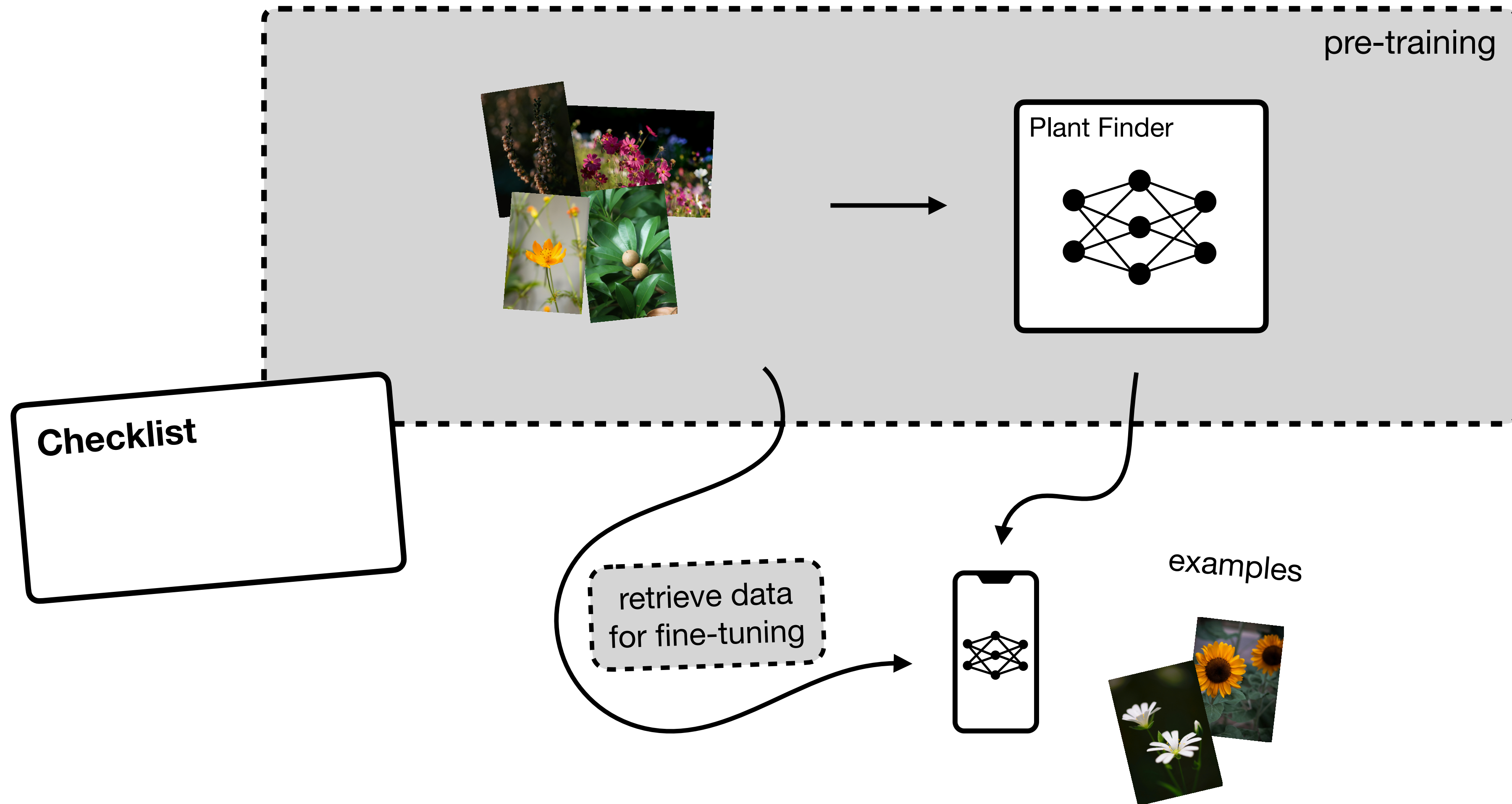
A Motivating Example



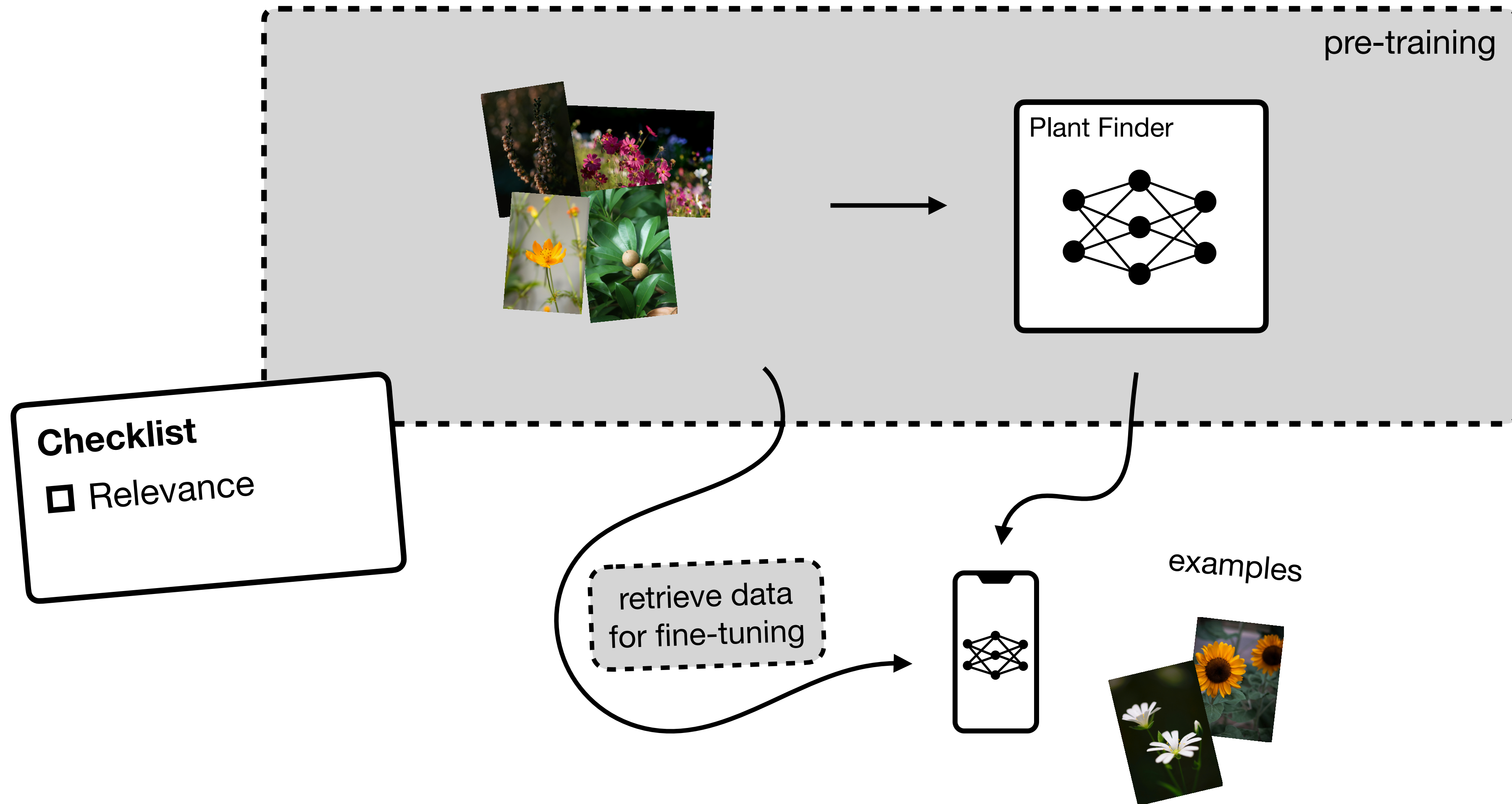
A Motivating Example



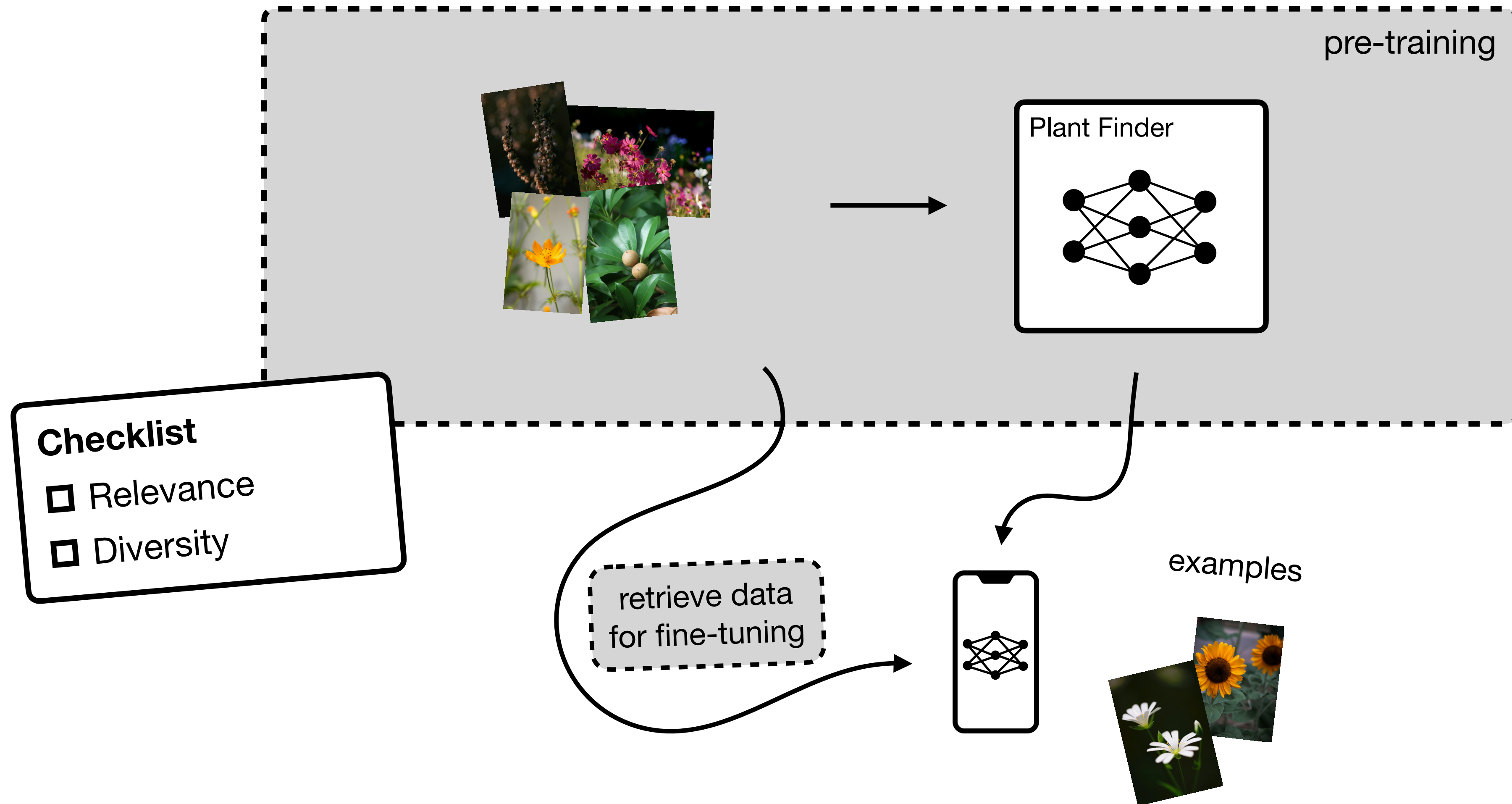
A Motivating Example



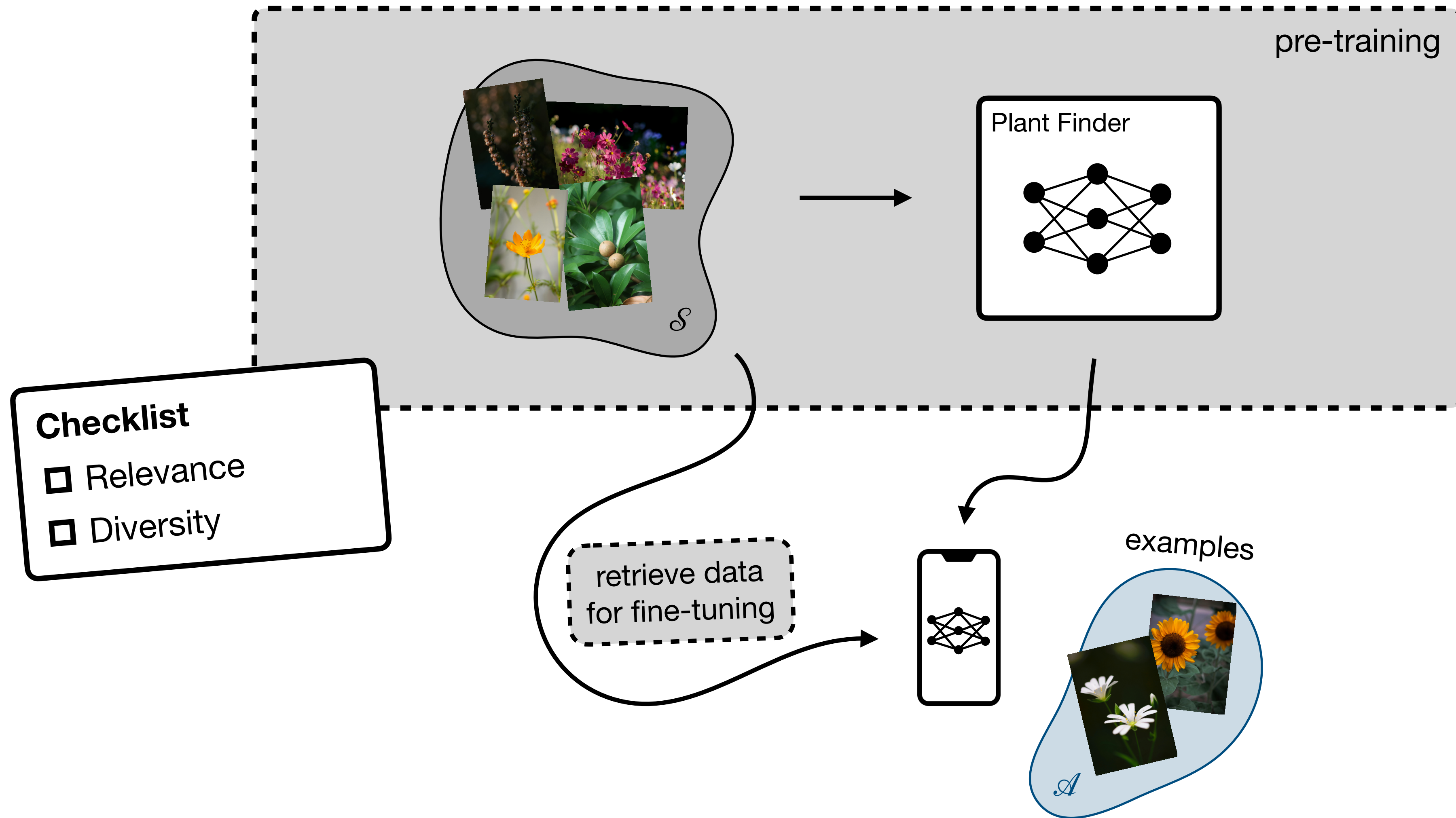
A Motivating Example



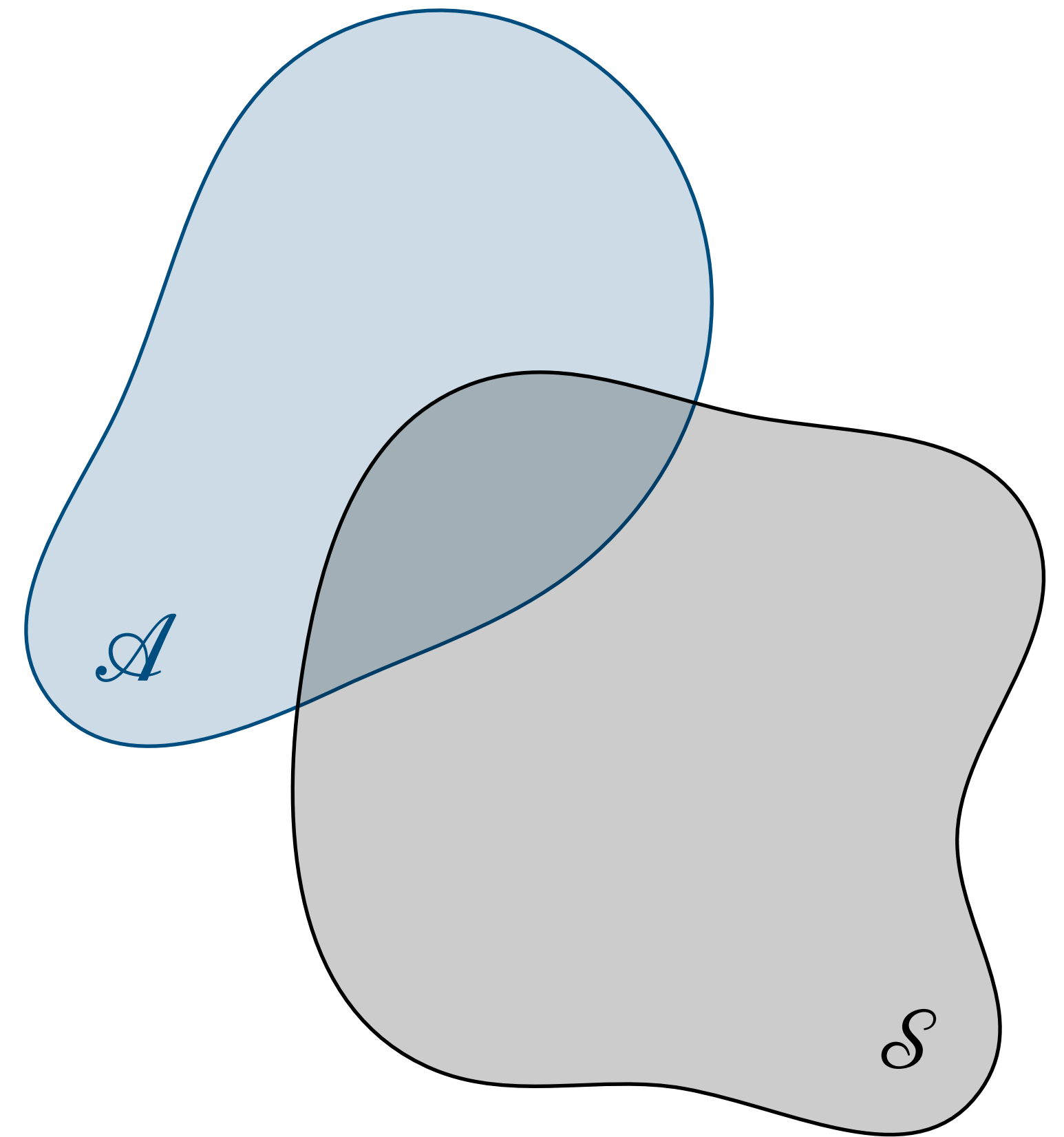
A Motivating Example



A Motivating Example

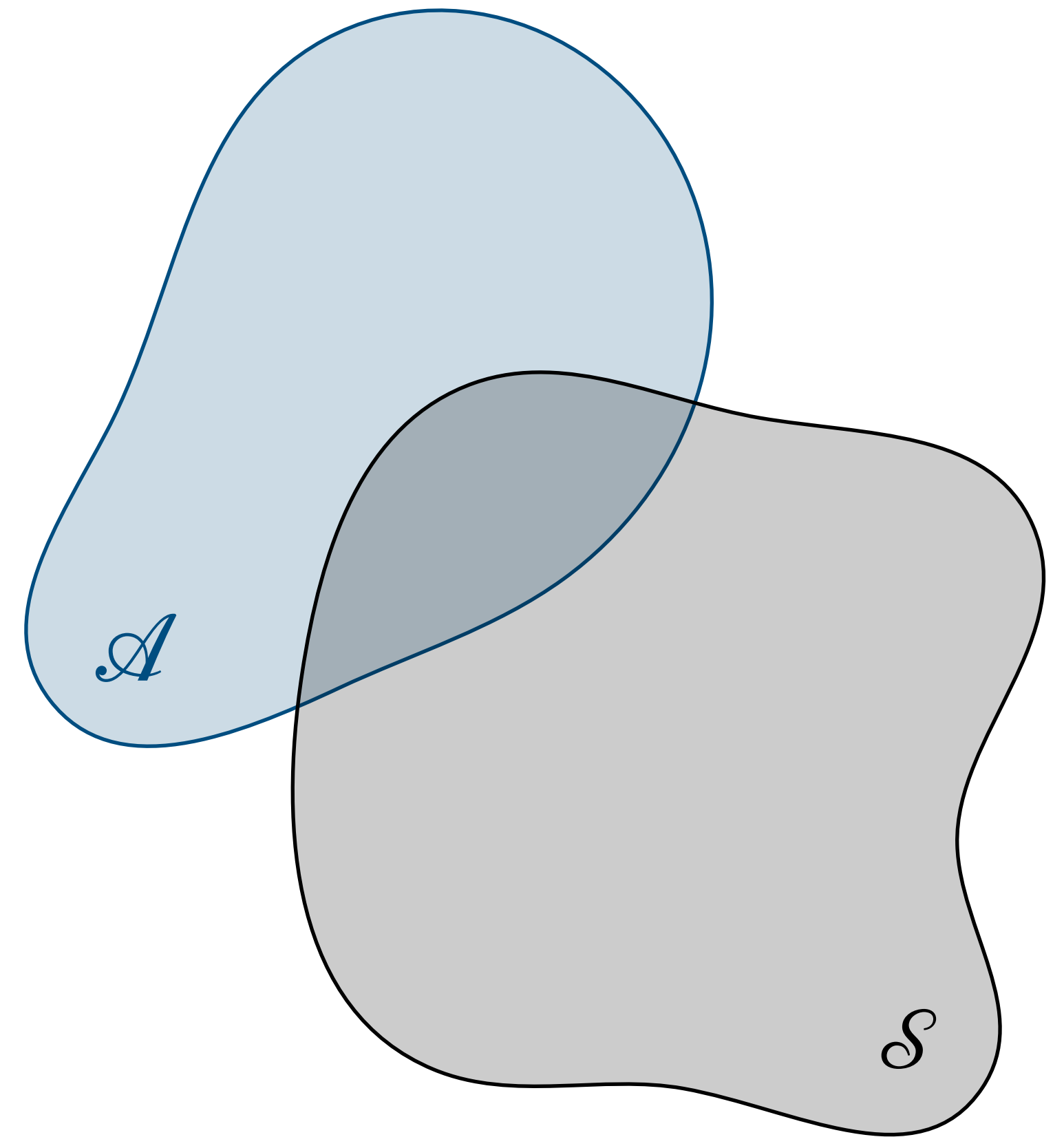


Setting



Setting

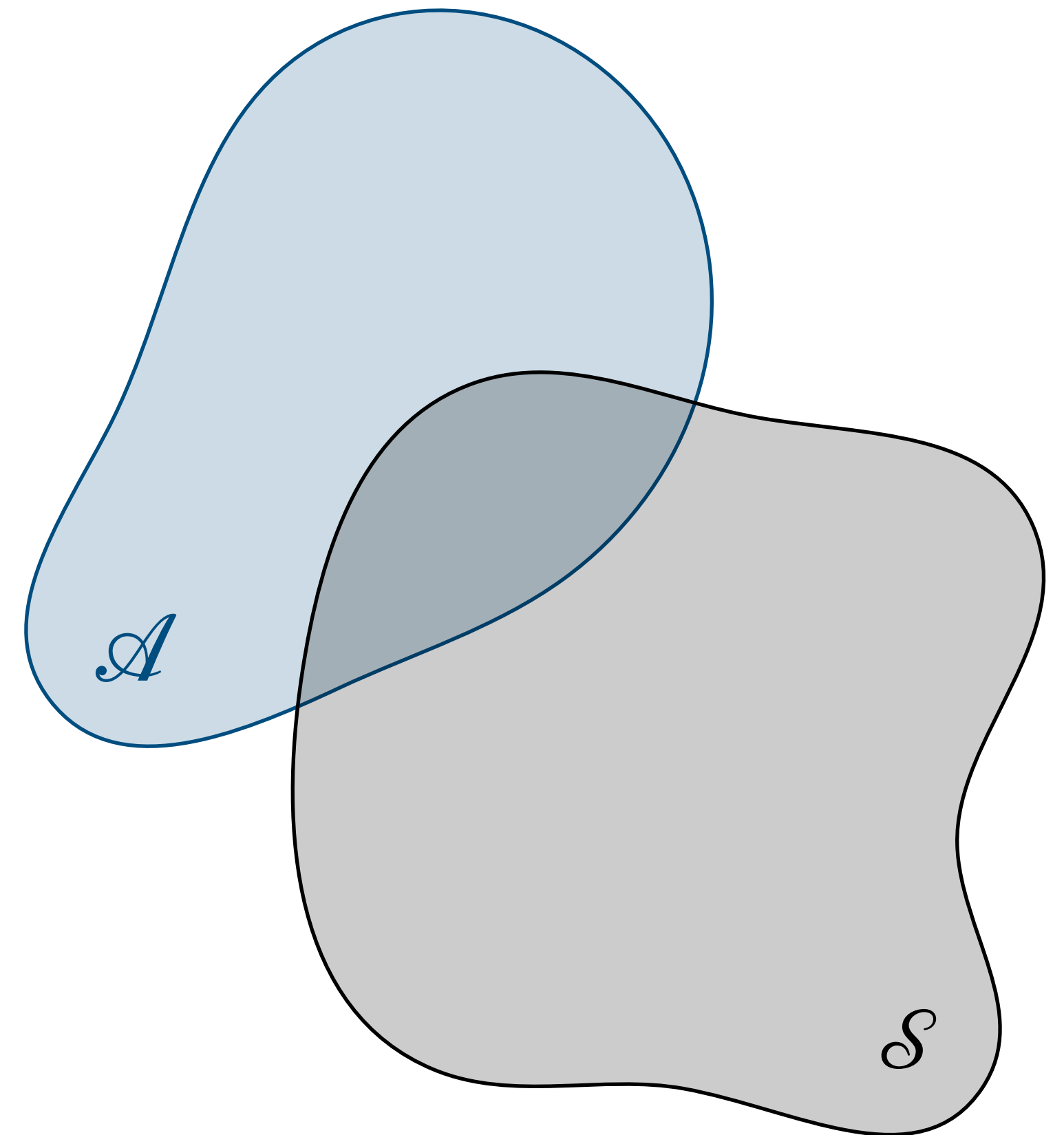
- $\mathcal{S} \subseteq \mathcal{X}$ - sample space
- $\mathcal{A} \subseteq \mathcal{X}$ - target space
- Unknown function f over \mathcal{X}



Setting

- $\mathcal{S} \subseteq \mathcal{X}$ - sample space
- $\mathcal{A} \subseteq \mathcal{X}$ - target space
- Unknown function f over \mathcal{X}

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

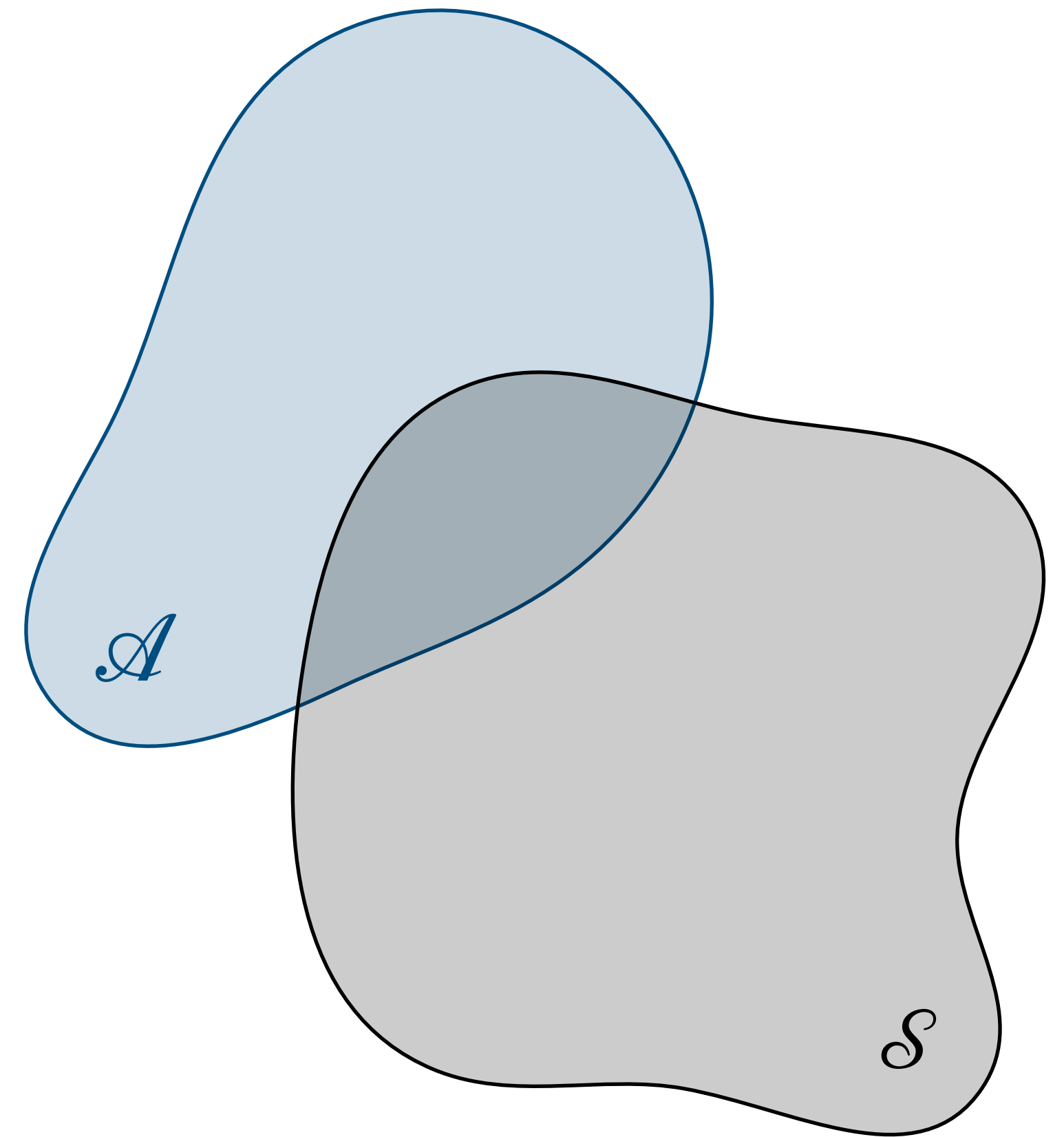


Setting

- $\mathcal{S} \subseteq \mathcal{X}$ - sample space
- $\mathcal{A} \subseteq \mathcal{X}$ - target space
- Unknown function f over \mathcal{X}

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

We call this **Transductive Active Learning**,
generalizing classical Active Learning



Setting

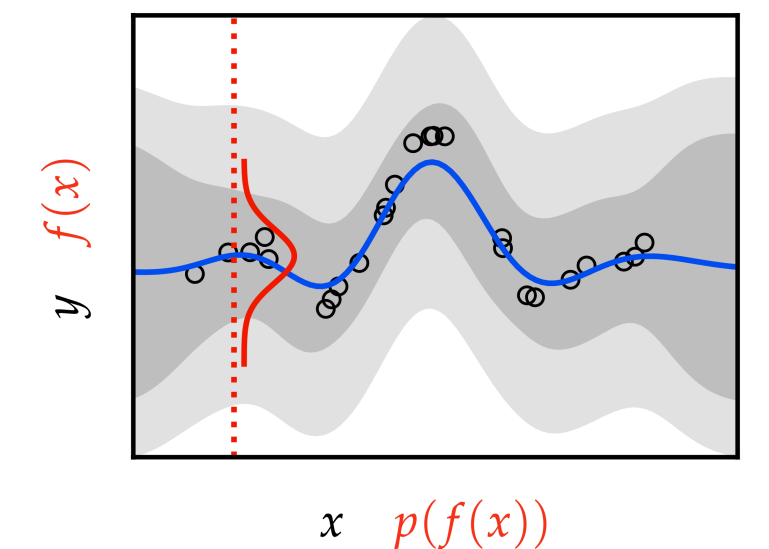
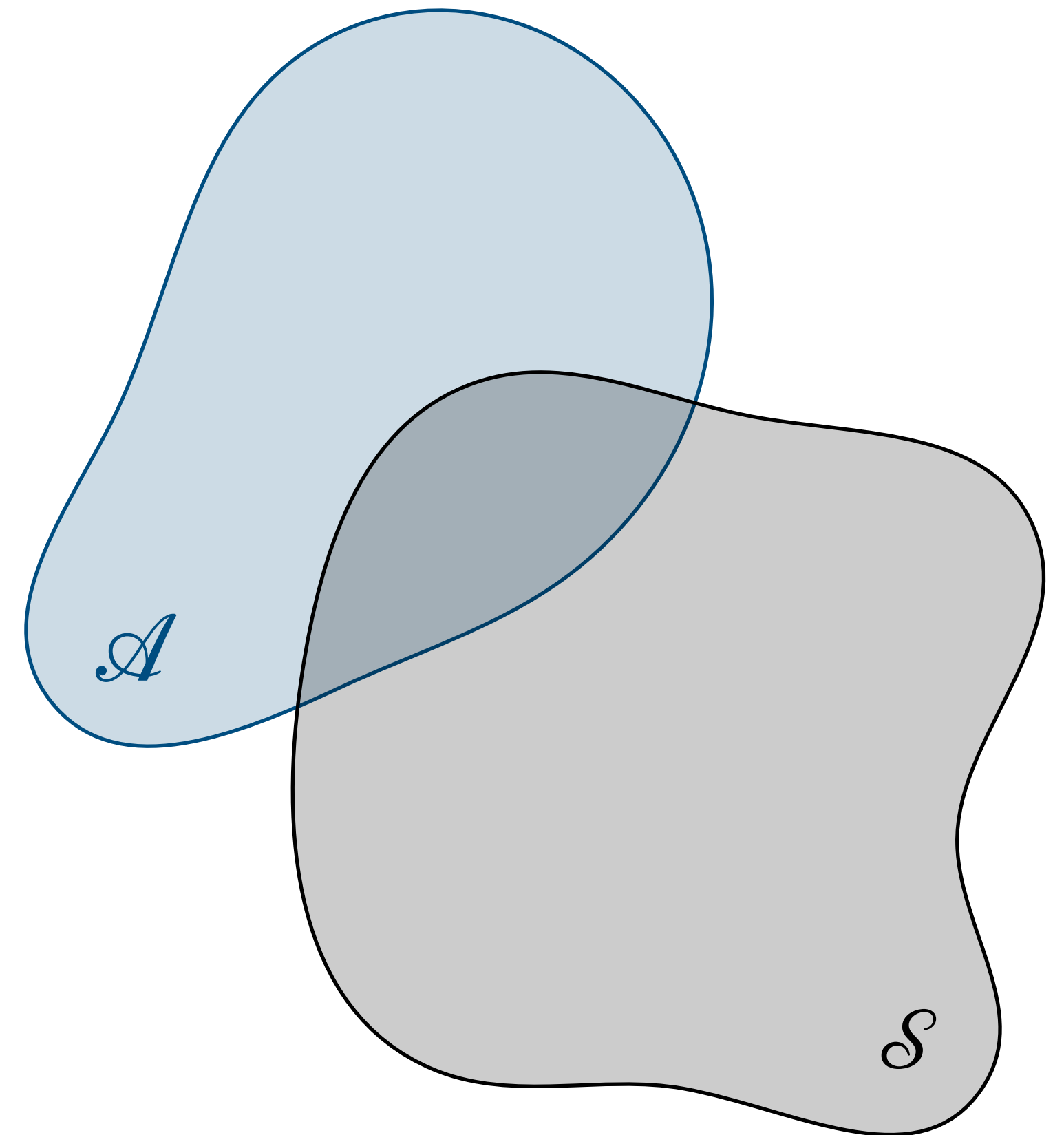
- $\mathcal{S} \subseteq \mathcal{X}$ - sample space
- $\mathcal{A} \subseteq \mathcal{X}$ - target space
- Unknown function f over \mathcal{X}

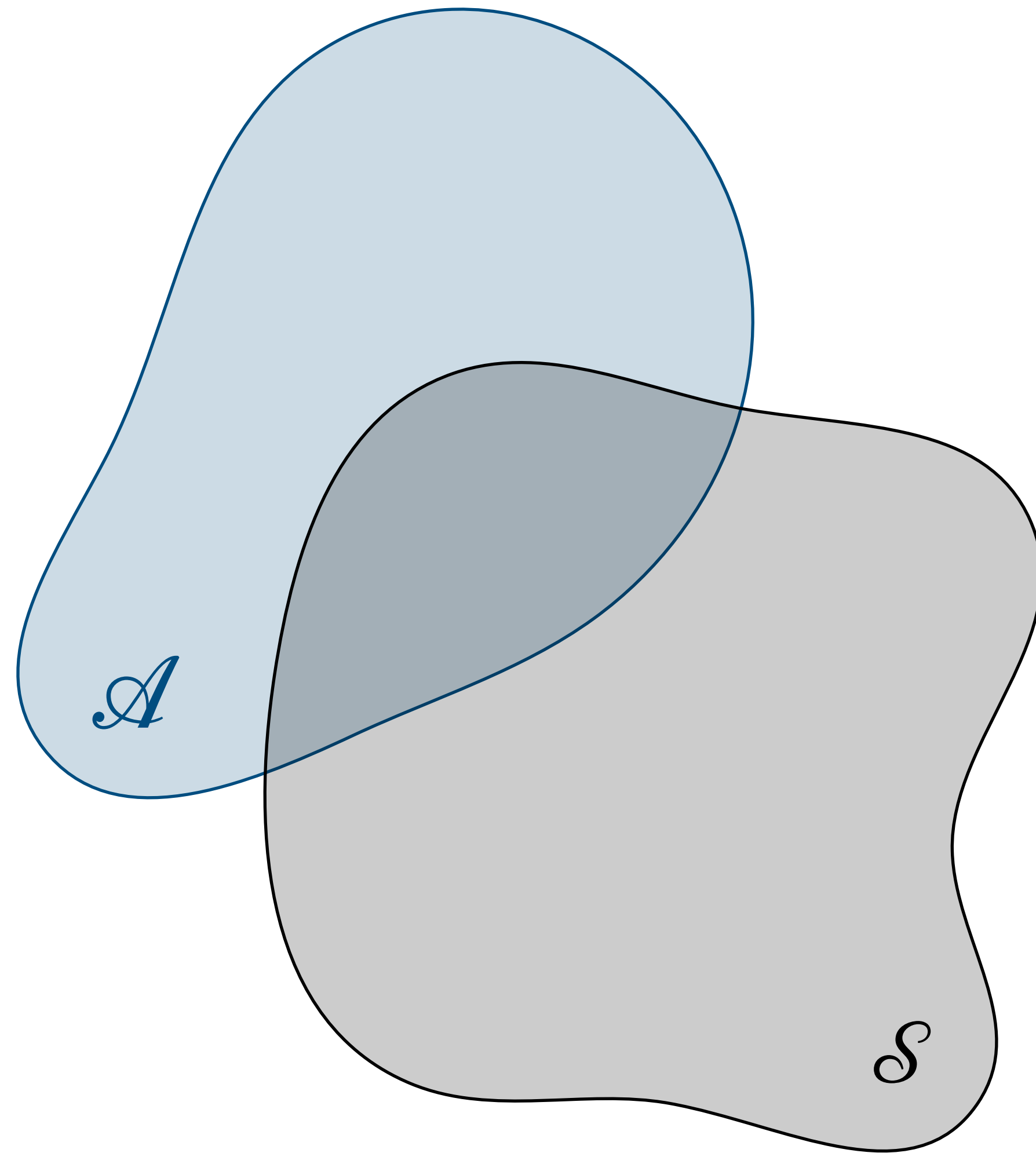
Goal: Learn f within \mathcal{A} by sampling from \mathcal{S}

We call this **Transductive Active Learning**, generalizing classical Active Learning

Assume for us: \mathcal{S}, \mathcal{A} finite, and NN f is approximated by a Gaussian process

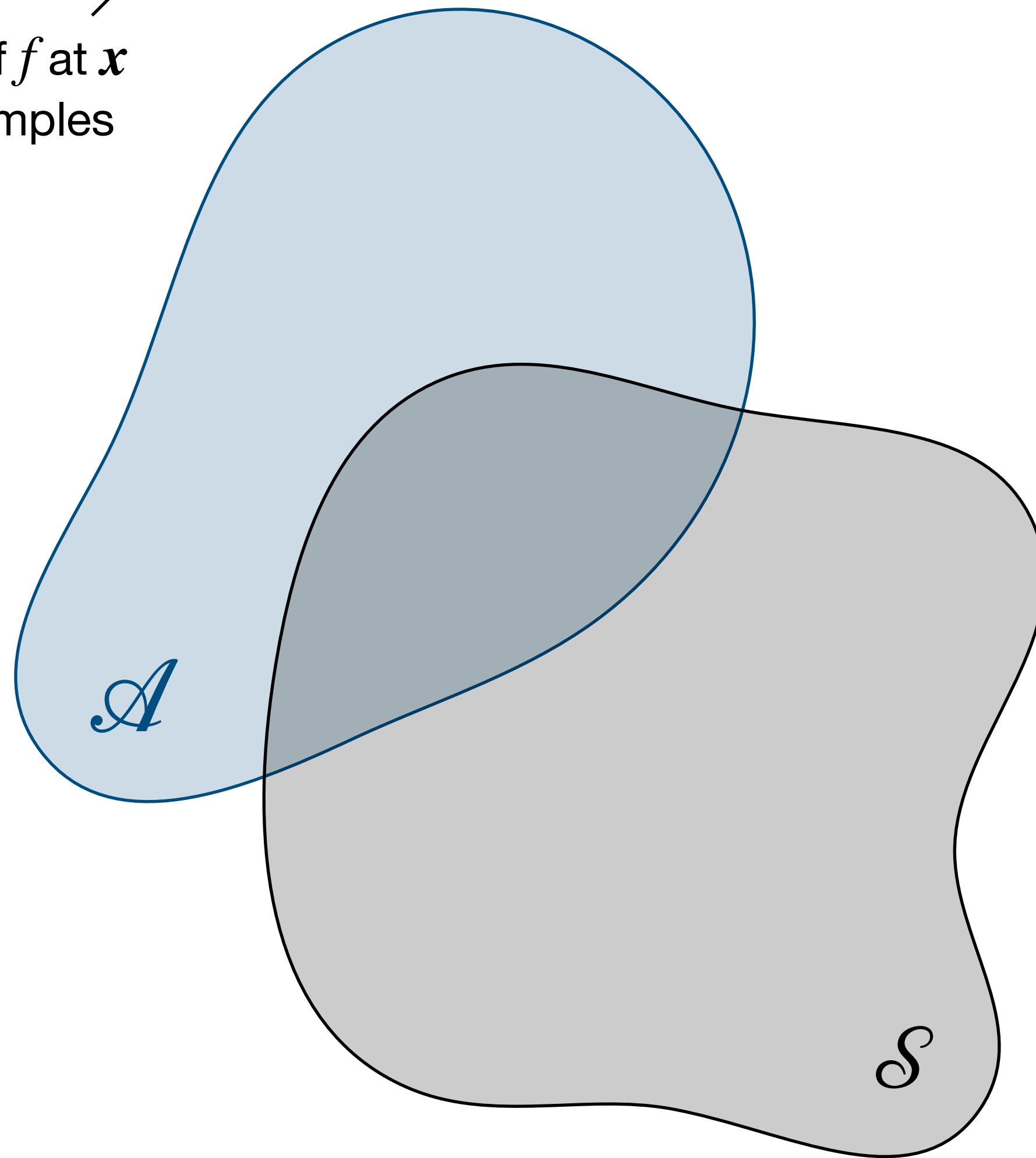
with kernel $k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}(\mathbf{x})^\top \boldsymbol{\phi}(\mathbf{x}')$ where $\boldsymbol{\phi}(\cdot)$ are embeddings generated by the NN





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

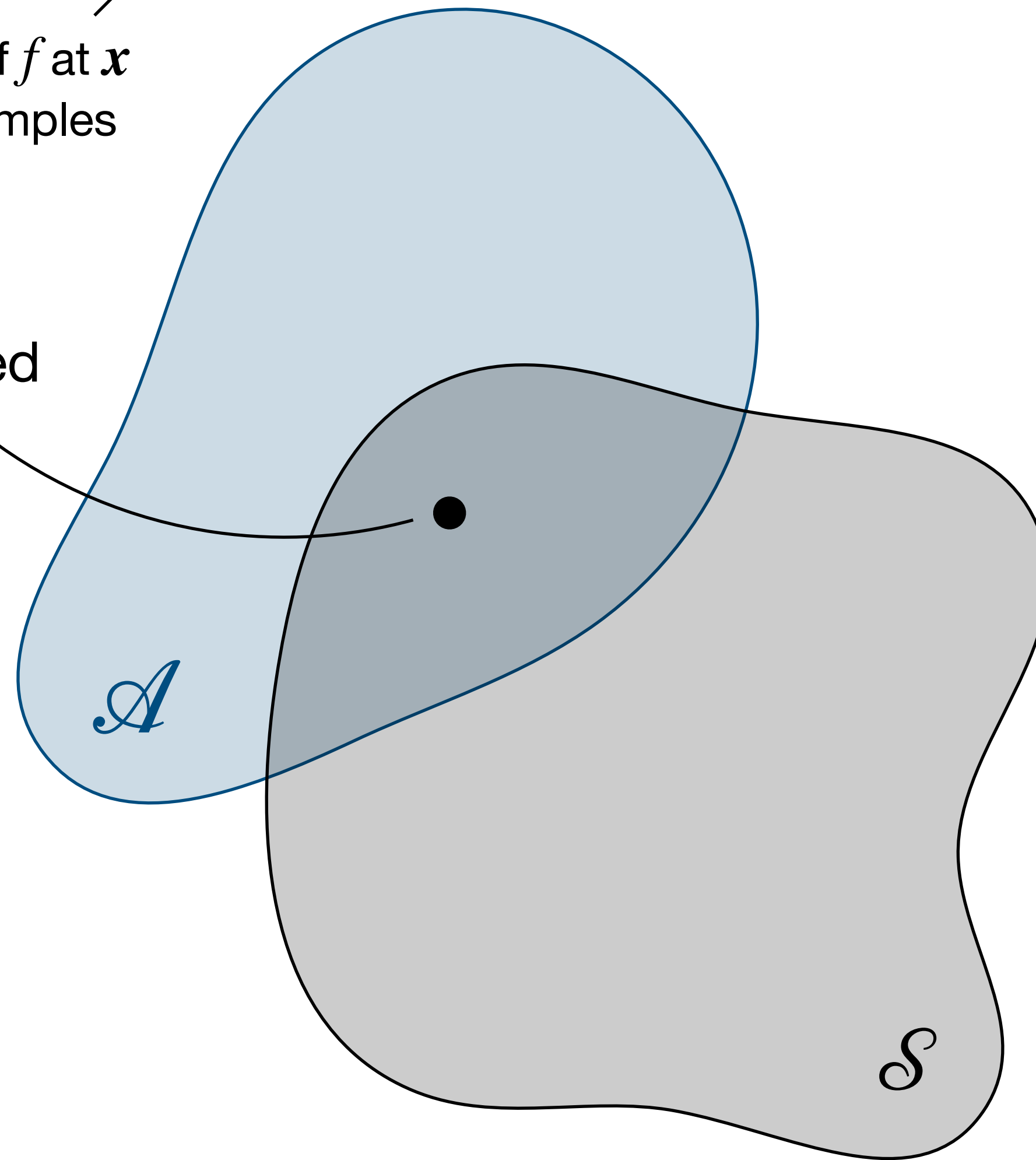
variance of f at \mathbf{x}
after n samples



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

variance of f at \mathbf{x}
after n samples

how much can be learned
about this point?

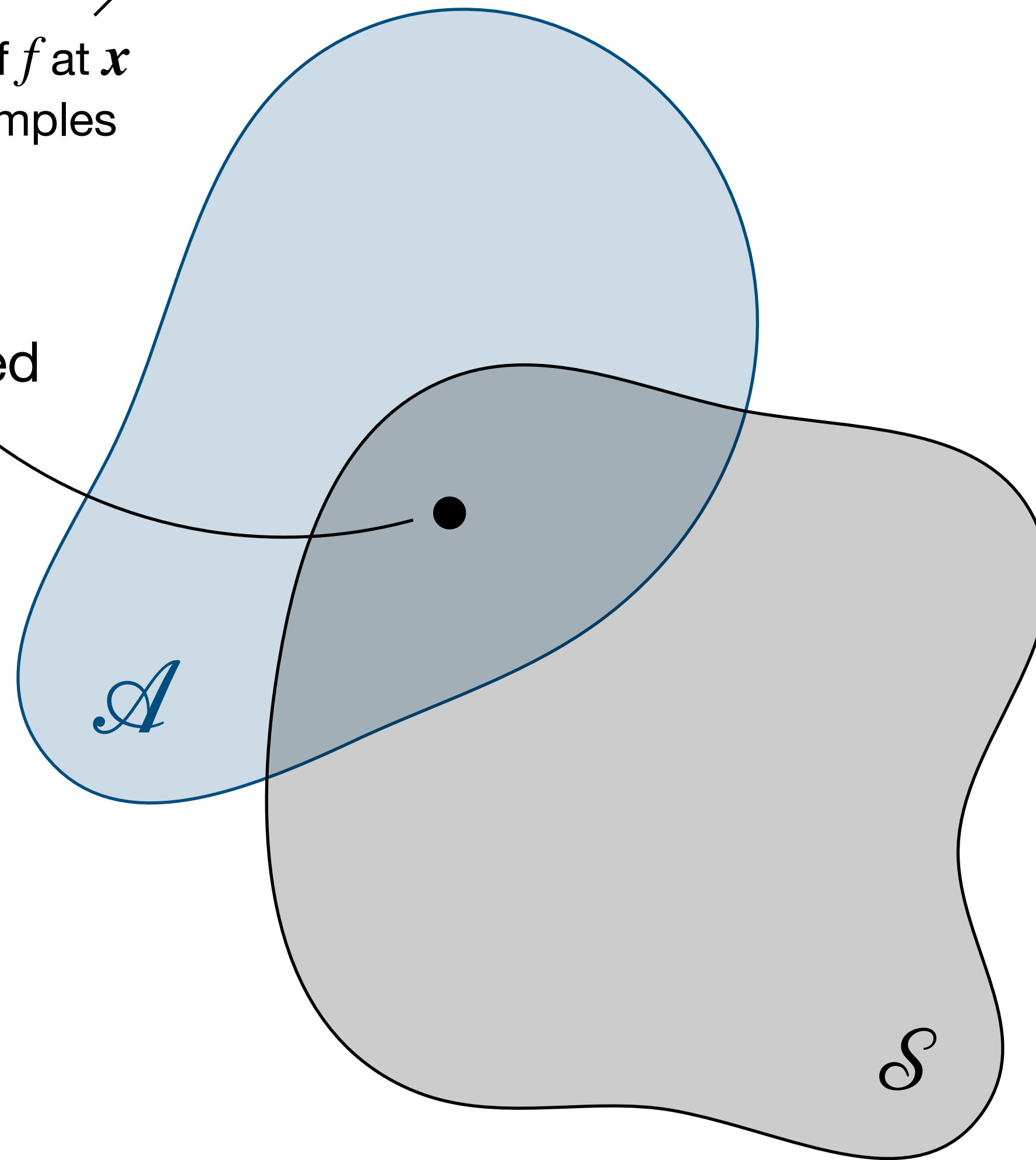


Goal: Reduce uncertainty $\sigma_n^2(\mathbf{x})$ at $\mathbf{x} \in \mathcal{A}$

variance of f at \mathbf{x}
after n samples

how much can be learned
about this point?

$\sigma_n^2(\mathbf{x}) \rightarrow 0$ as $n \rightarrow \infty$

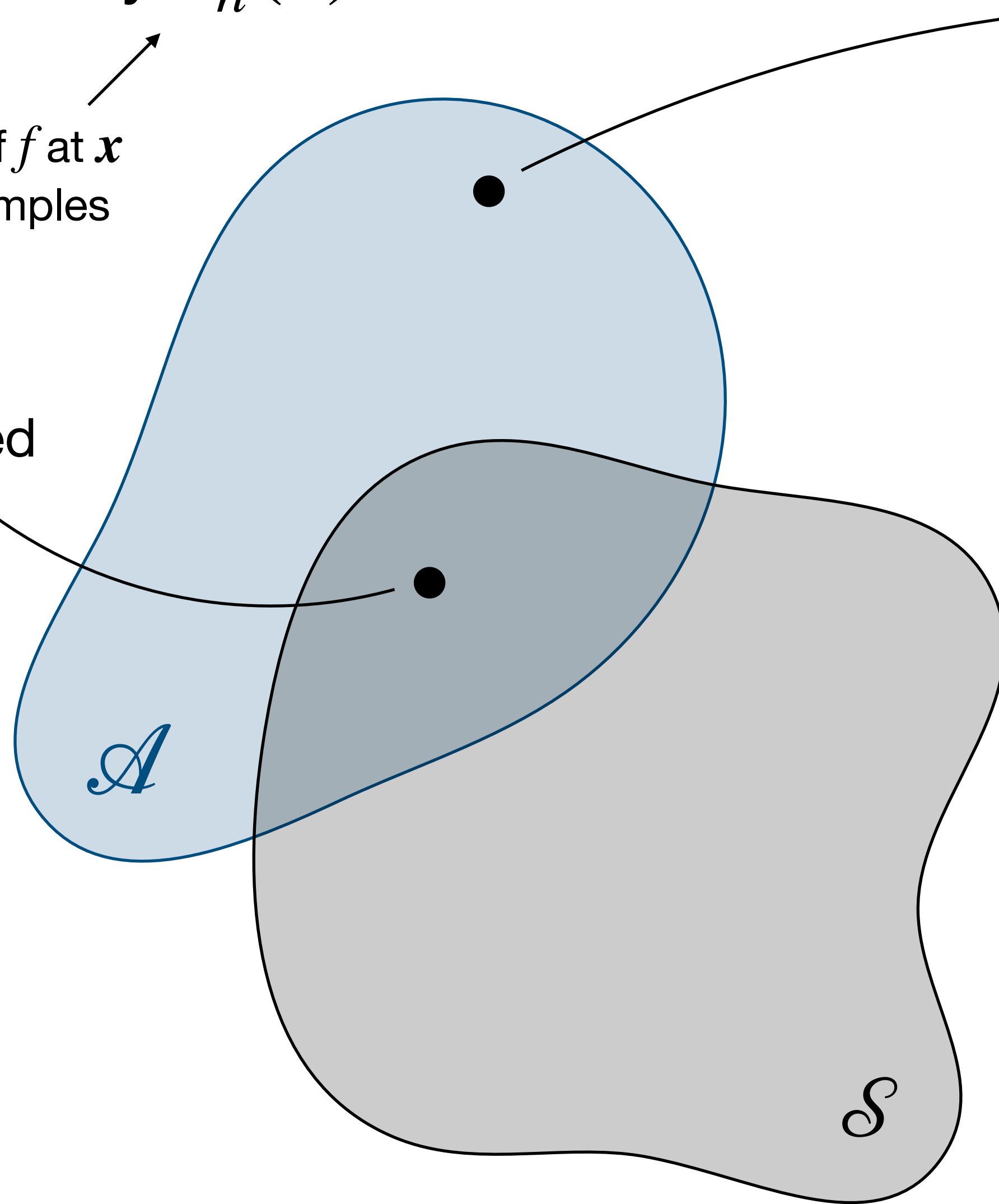


Goal: Reduce uncertainty $\sigma_n^2(\mathbf{x})$ at $\mathbf{x} \in \mathcal{A}$

variance of f at \mathbf{x}
after n samples

how much can be learned
about this point?

$\sigma_n^2(\mathbf{x}) \rightarrow 0$ as $n \rightarrow \infty$



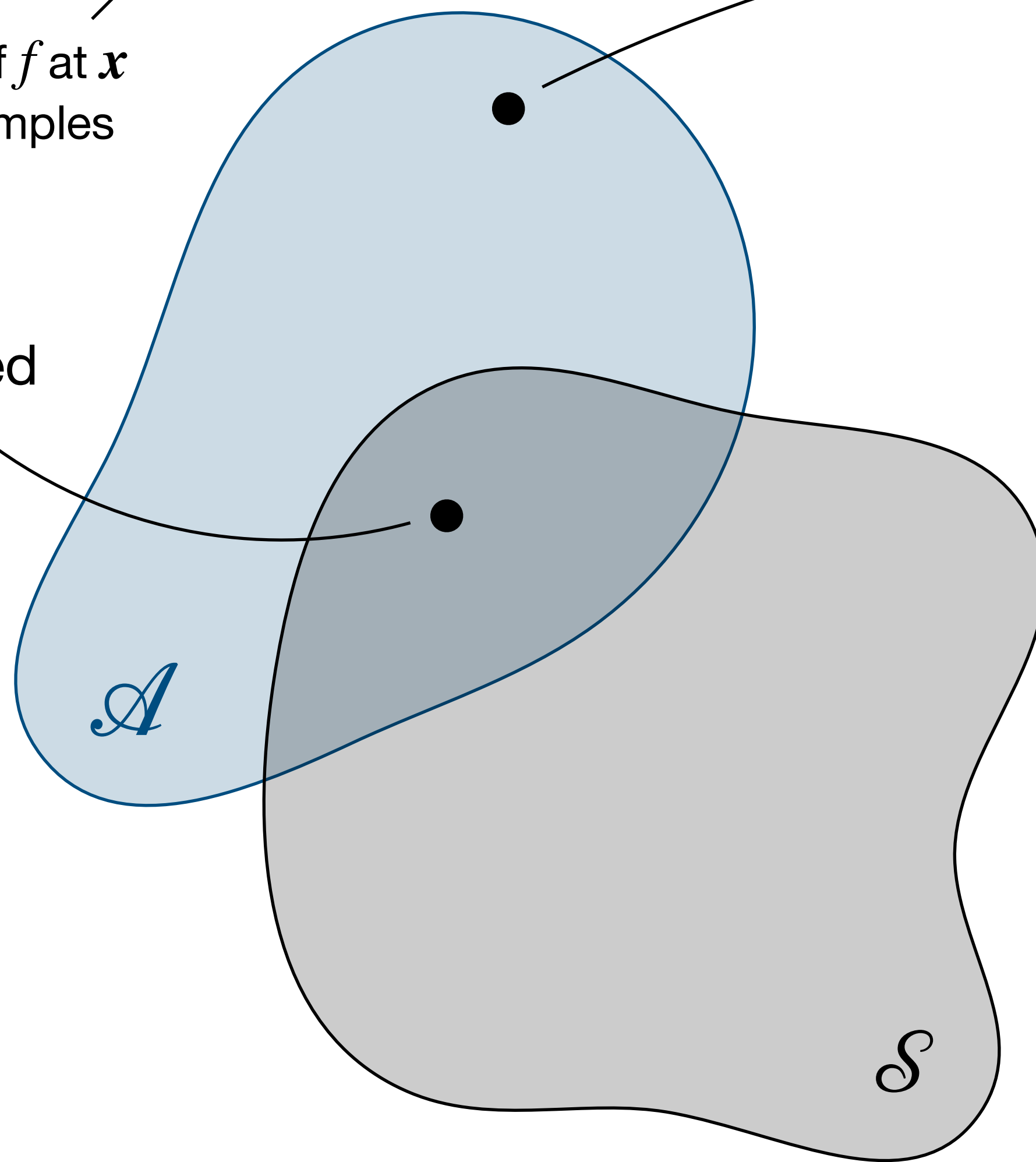
what about the point \mathbf{x}' ?

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

variance of f at \mathbf{x}
after n samples

how much can be learned
about this point?

$\sigma_n^2(\mathbf{x}) \rightarrow 0$ as $n \rightarrow \infty$



what about the point \mathbf{x}' ?

$\sigma_n^2(\mathbf{x}') \xrightarrow{?} \eta_{\mathcal{S}}^2(\mathbf{x}') \text{ as } n \rightarrow \infty$

where $\eta_{\mathcal{S}}^2(\mathbf{x}') = \text{Var}[f(\mathbf{x}') \mid f(\mathcal{S})]$ is
the **irreducible uncertainty**

ITL: Information-directed Transductive Learning

Proposal: select the next sample to minimize *posterior* uncertainty within \mathcal{A}

ITL: Information-directed Transductive Learning

Proposal: select the next sample to minimize *posterior* uncertainty within \mathcal{A}

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

ITL: Information-directed Transductive Learning

Proposal: select the next sample to minimize *posterior* uncertainty within \mathcal{A}

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

MacKay, 1992



ITL: Information-directed Transductive Learning

Proposal: select the next sample to minimize *posterior* uncertainty within \mathcal{A}

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

Generalization bound for ITL (informal). $\forall \mathbf{x}' \in \mathcal{A}$:

$$\sigma_n^2(\mathbf{x}') \leq \eta_{\mathcal{S}}^2(\mathbf{x}') + C \log n / \sqrt{n} \quad (C \text{ is a constant})$$

ITL: Information-directed Transductive Learning

Proposal: select the next sample to minimize *posterior* uncertainty within \mathcal{A}

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

Generalization bound for ITL (informal). $\forall \mathbf{x}' \in \mathcal{A}$:

$$\sigma_n^2(\mathbf{x}') \leq \underbrace{\eta_{\mathcal{S}}^2(\mathbf{x}')}_{\text{irreducible}} + \underbrace{C \log n / \sqrt{n}}_{\text{reducible}} \quad (C \text{ is a constant})$$

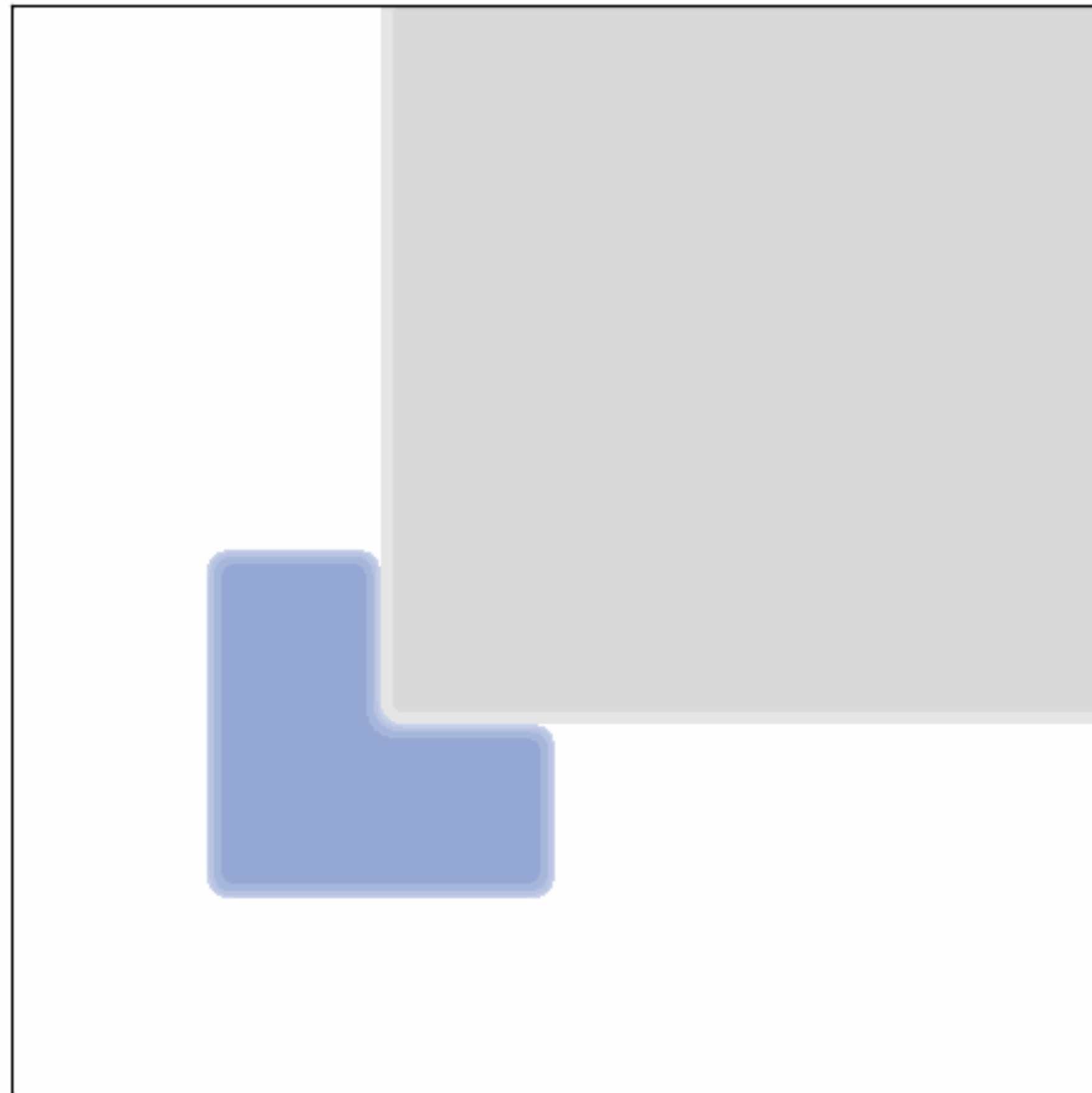
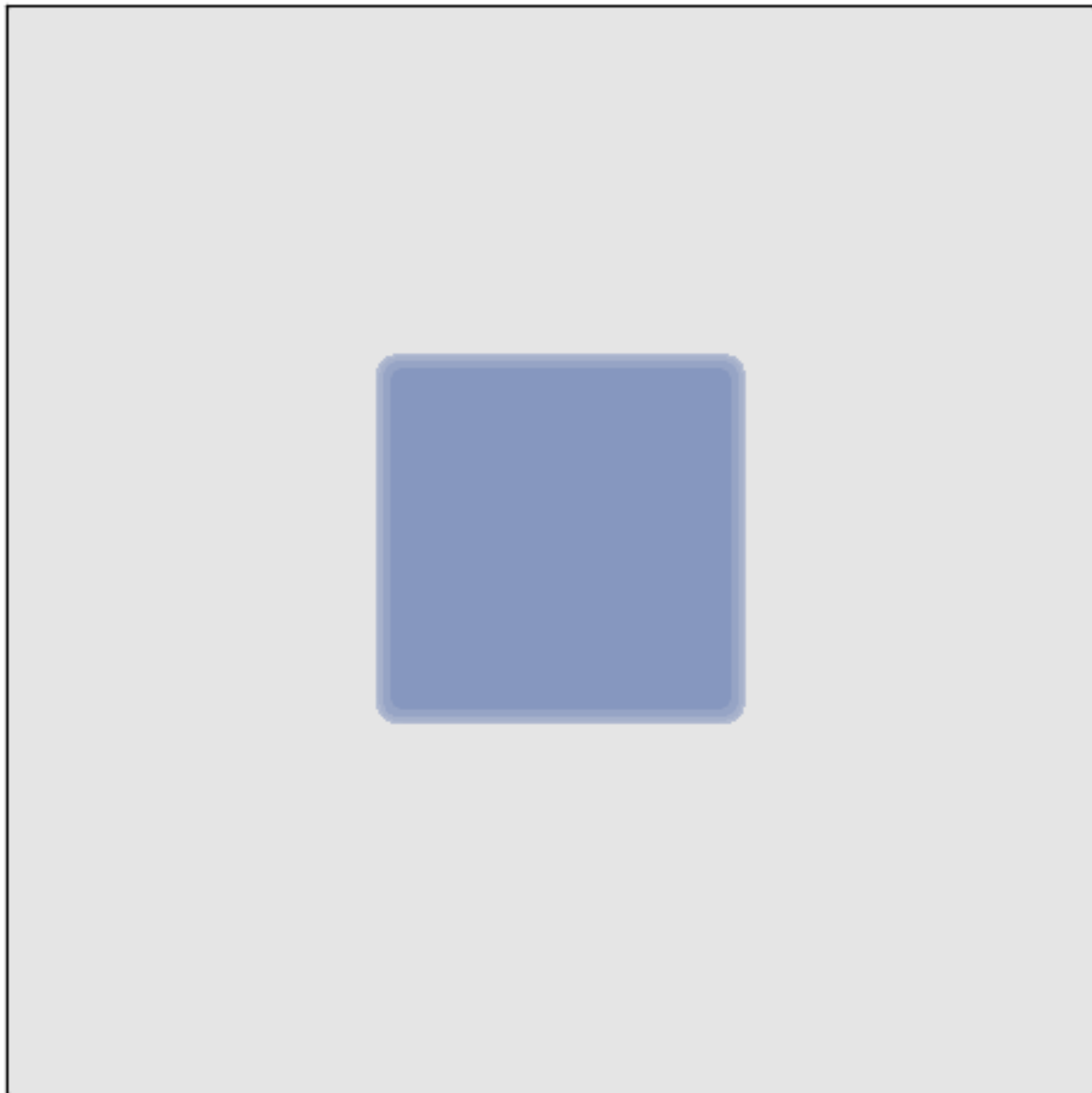


Illustration of ITL on a Gaussian process with Gaussian kernel

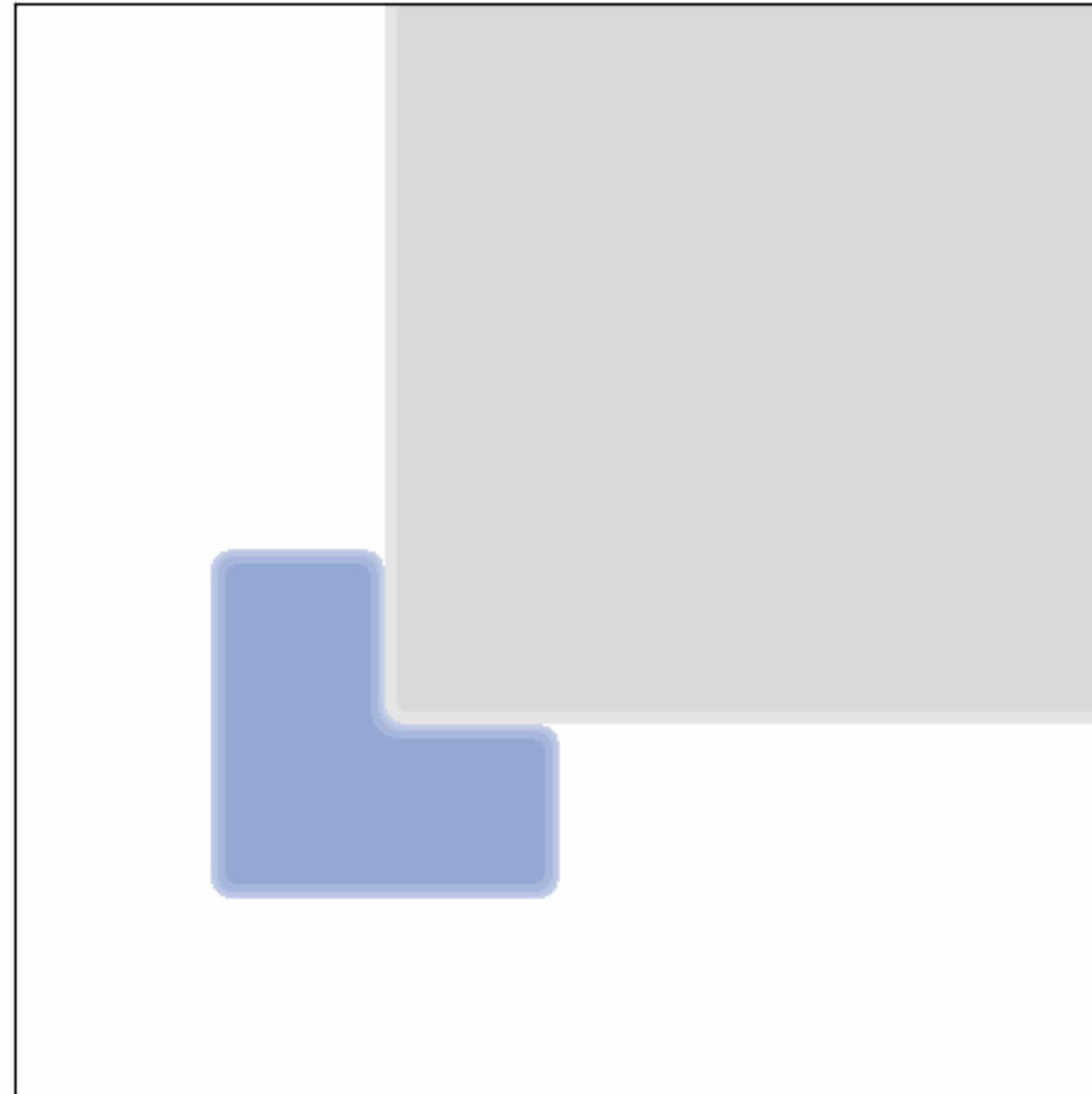
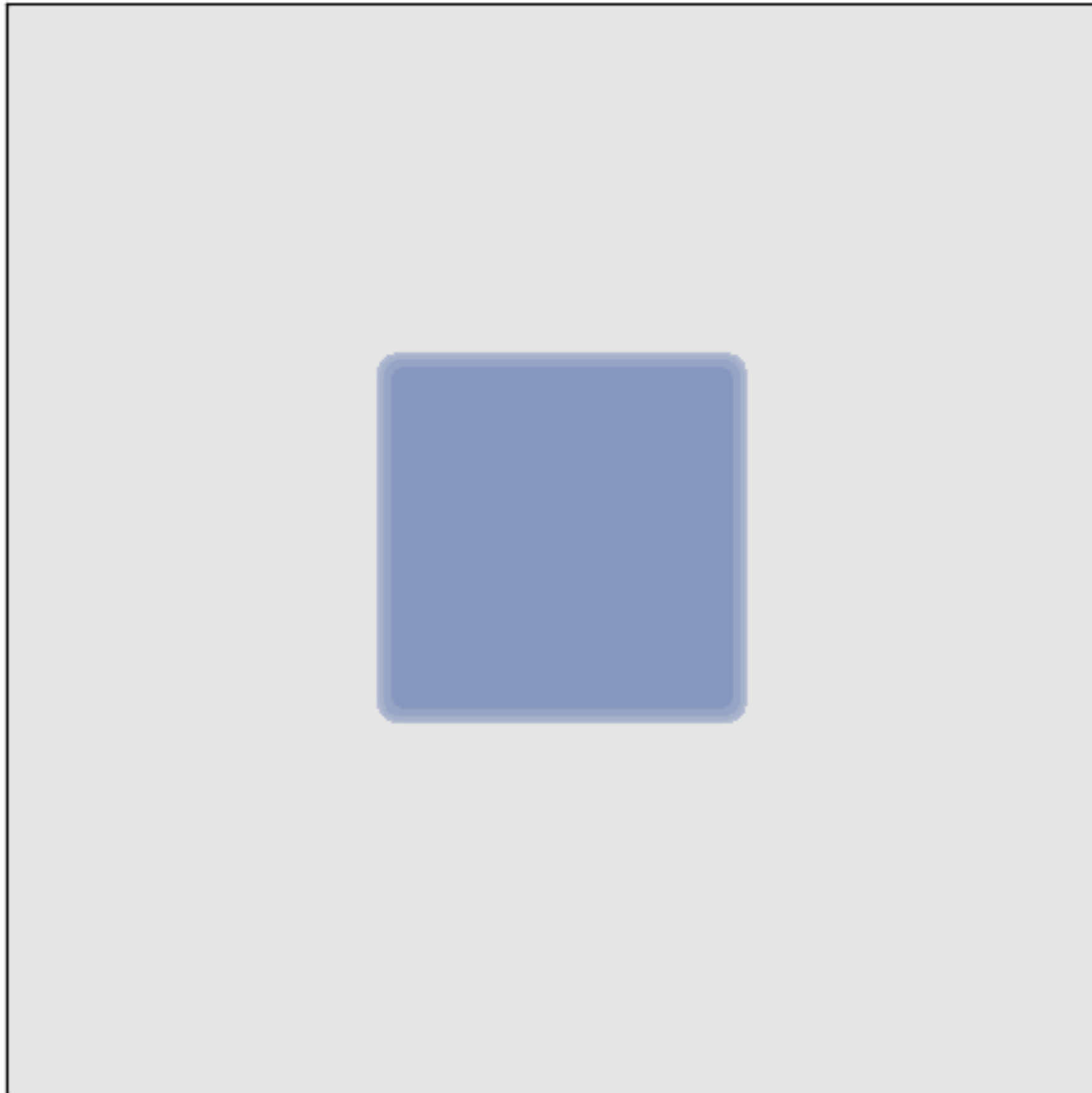
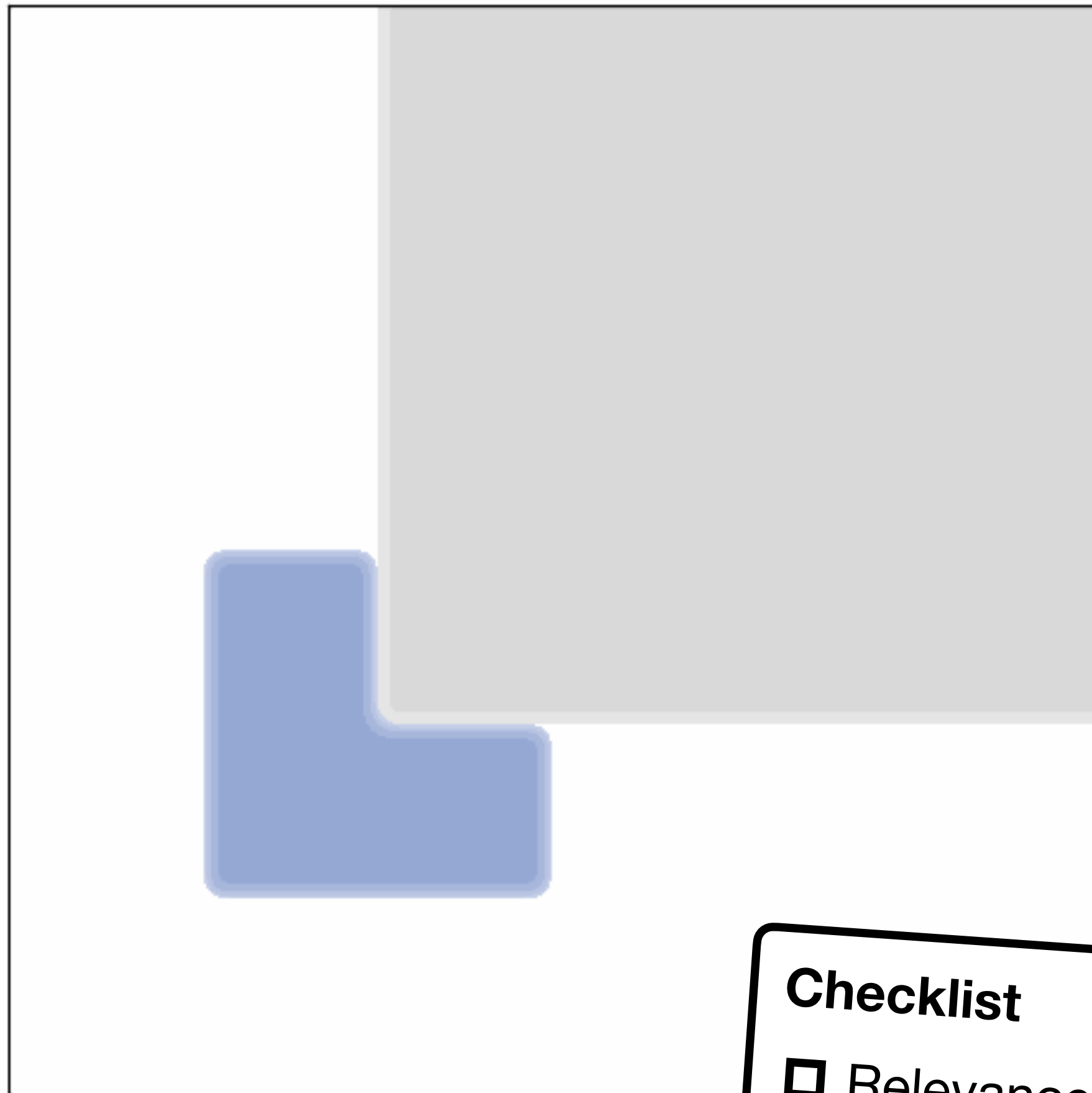
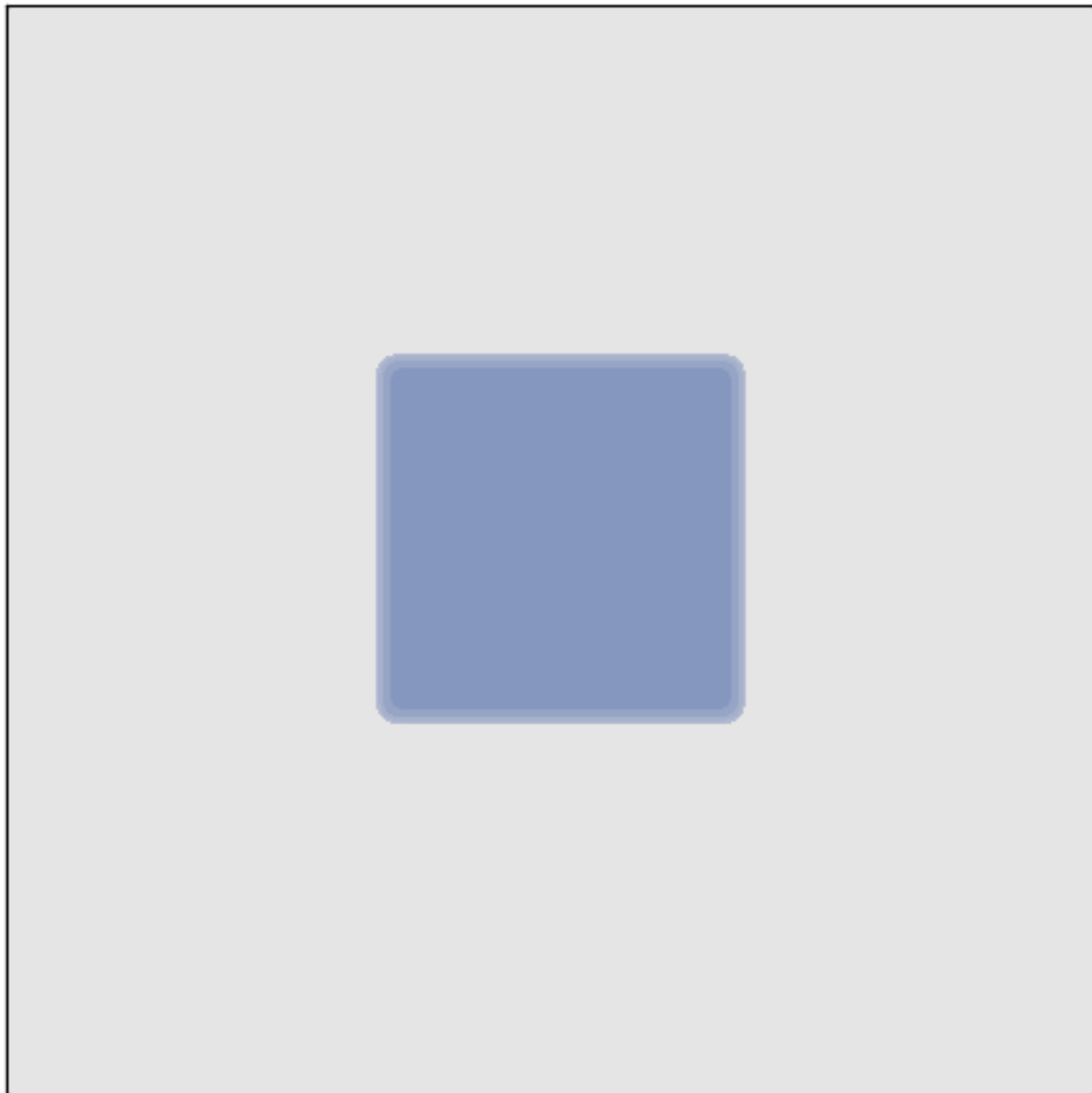
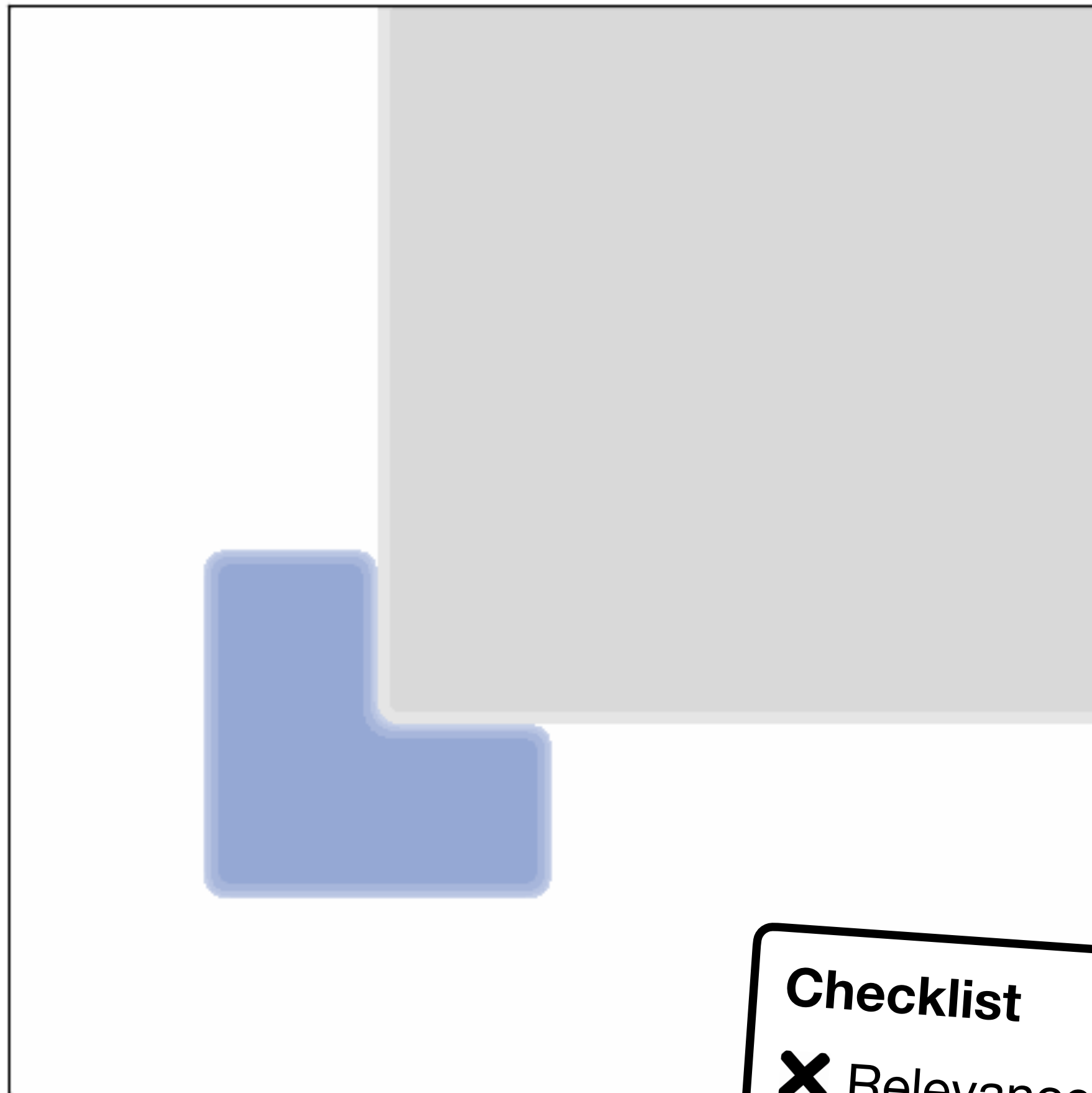
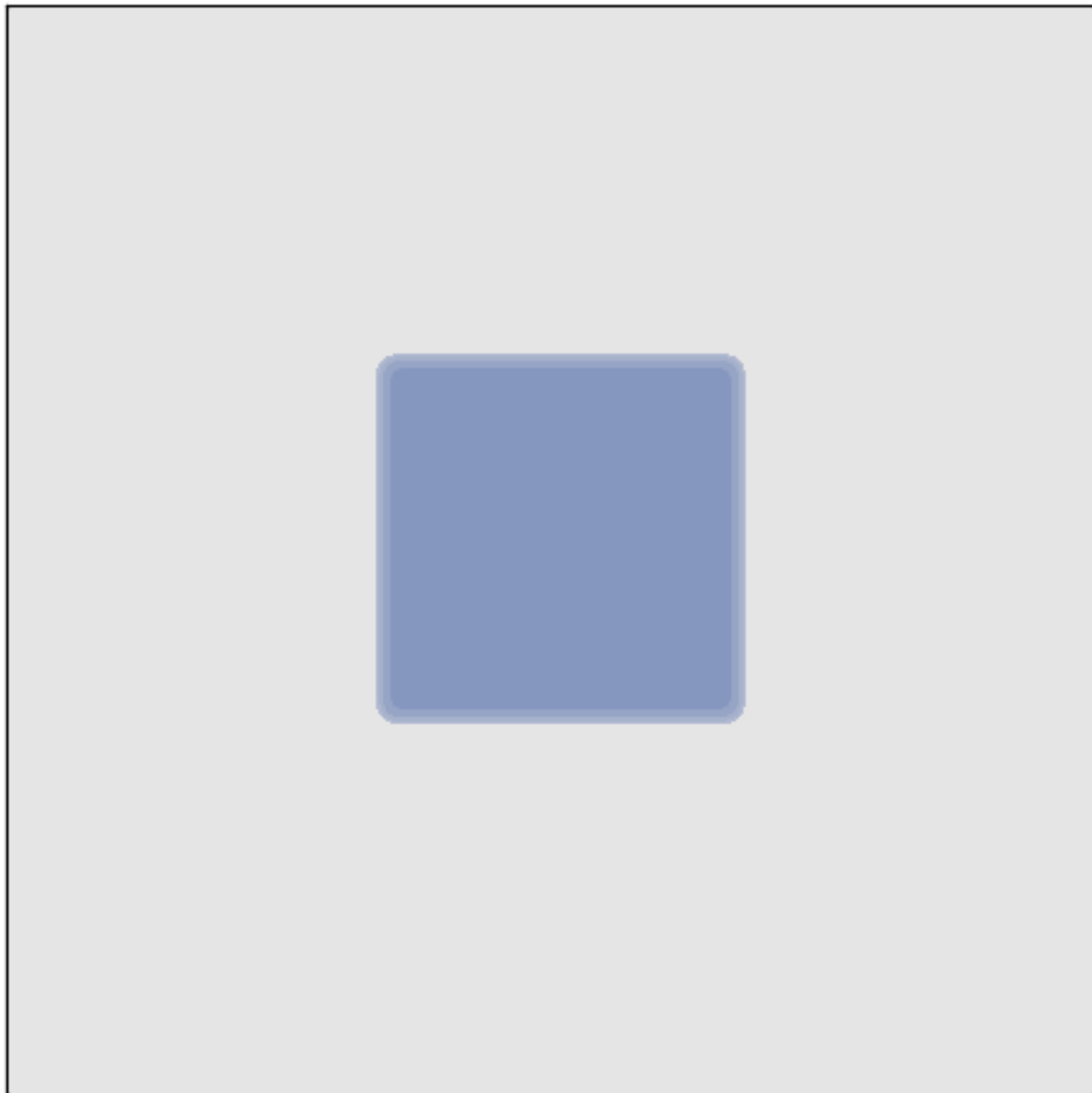


Illustration of ITL on a Gaussian process with Gaussian kernel



- Checklist**
- Relevance
 - Diversity

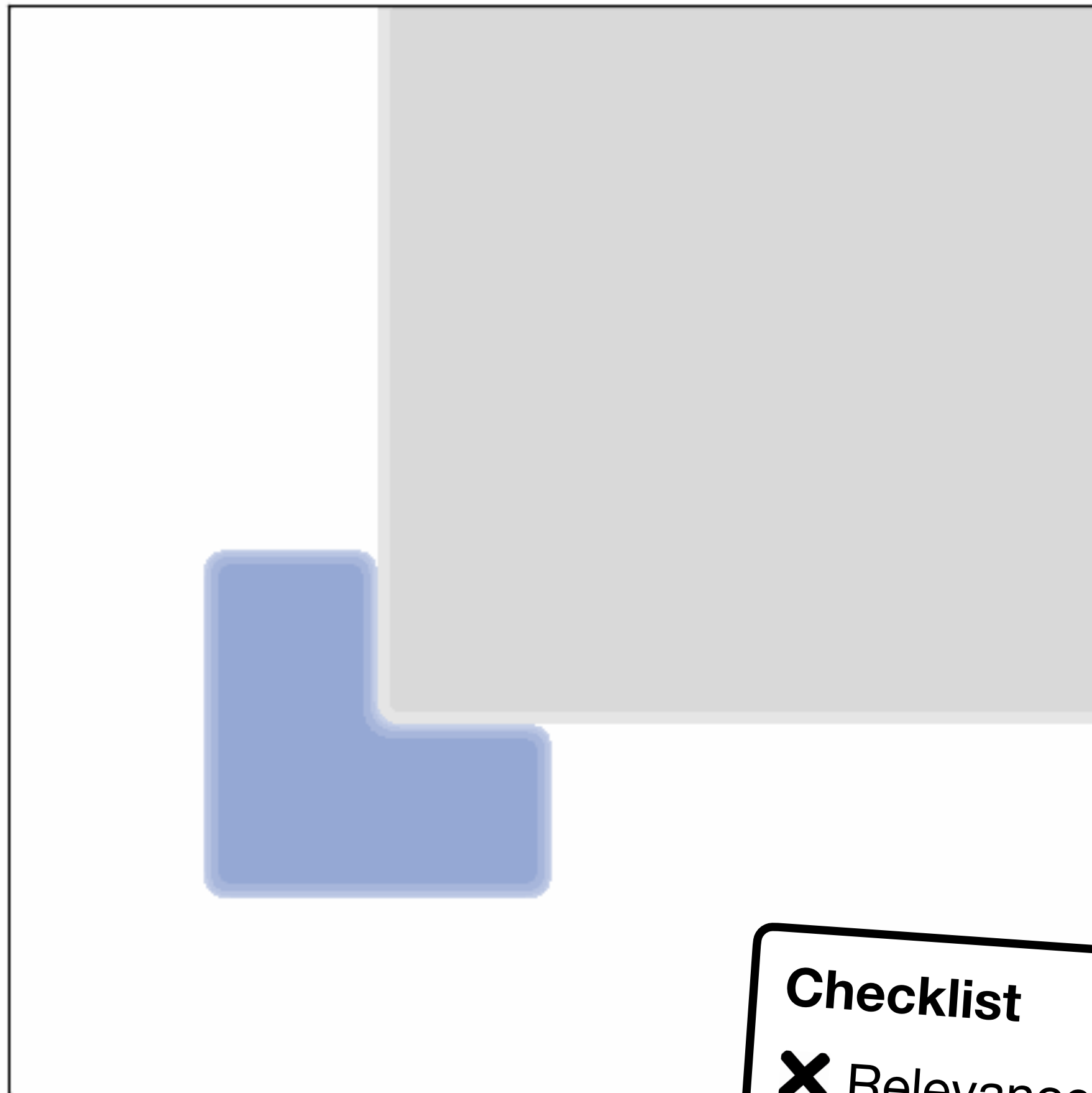
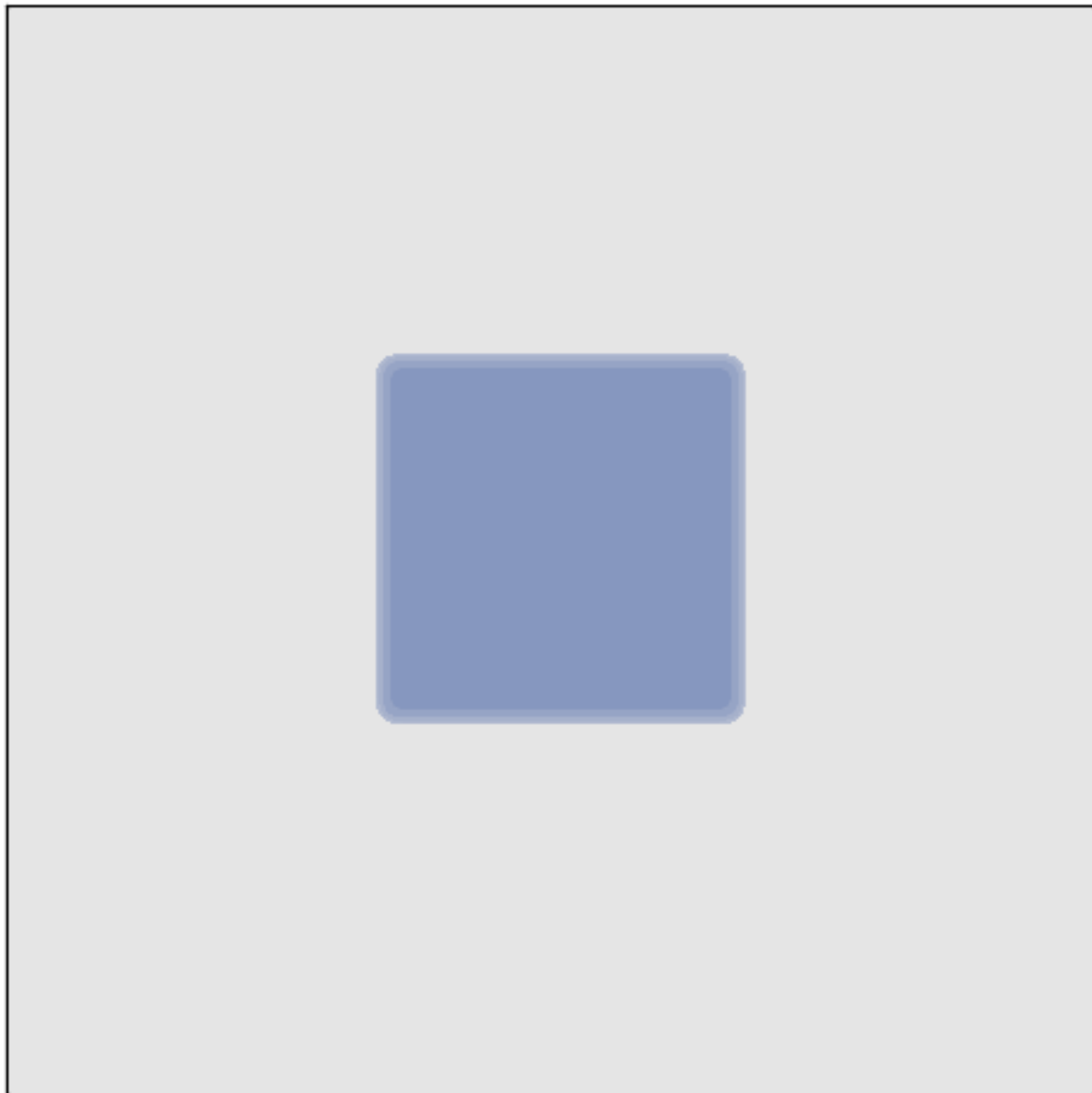
Illustration of ITL on a Gaussian process with Gaussian kernel



Checklist

- Relevance
- Diversity

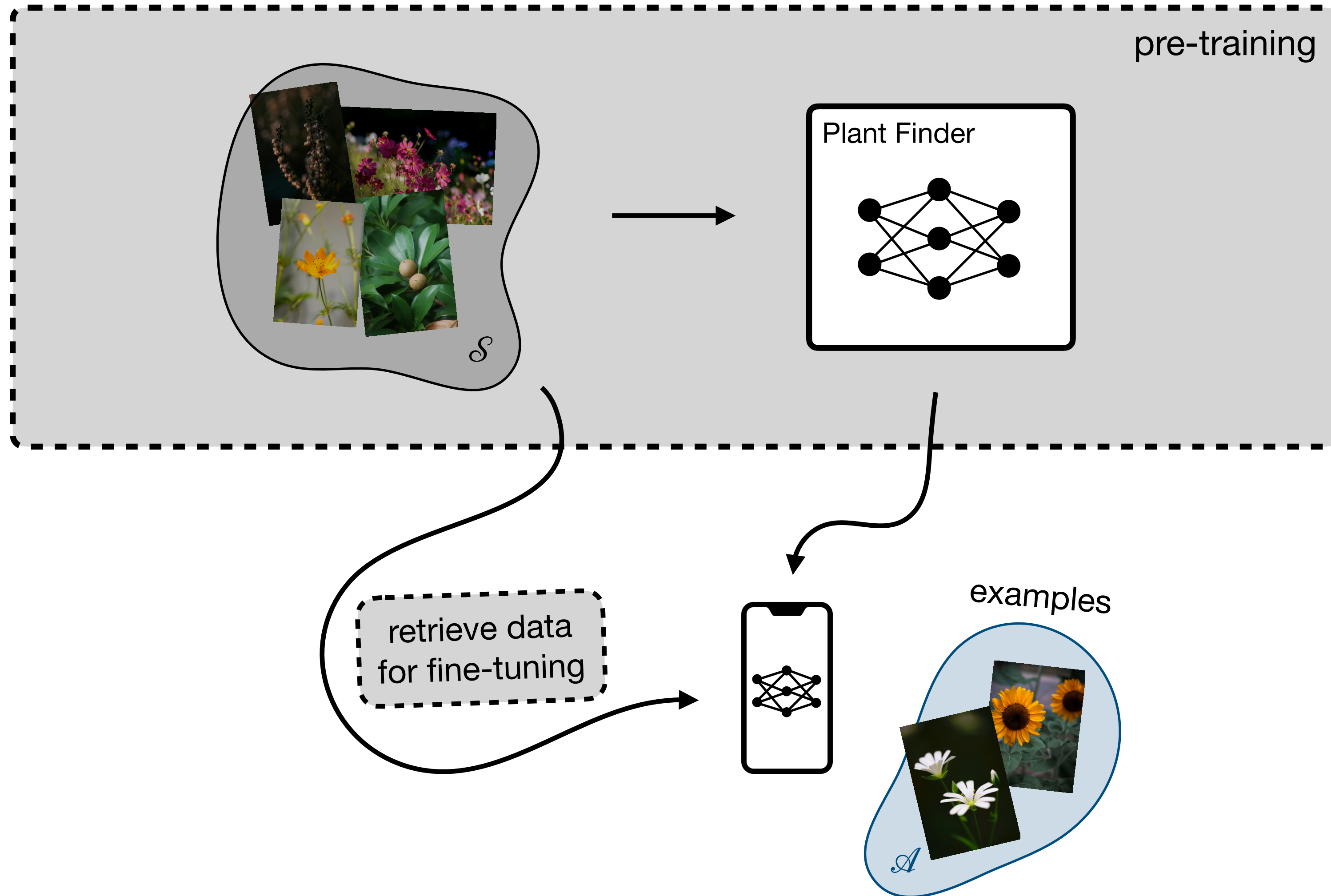
Illustration of ITL on a Gaussian process with Gaussian kernel



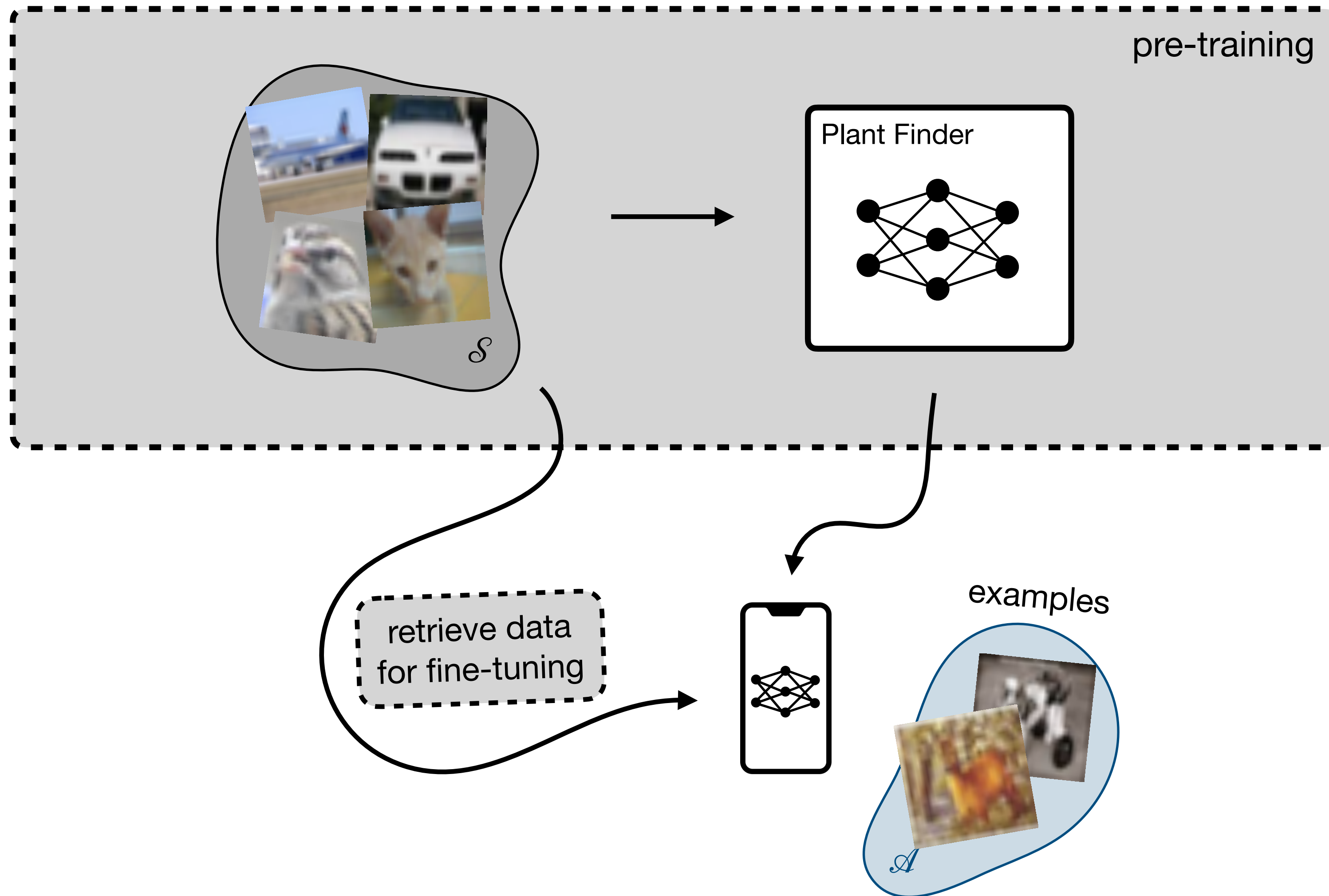
- Checklist**
- Relevance
 - Diversity

Illustration of ITL on a Gaussian process with Gaussian kernel

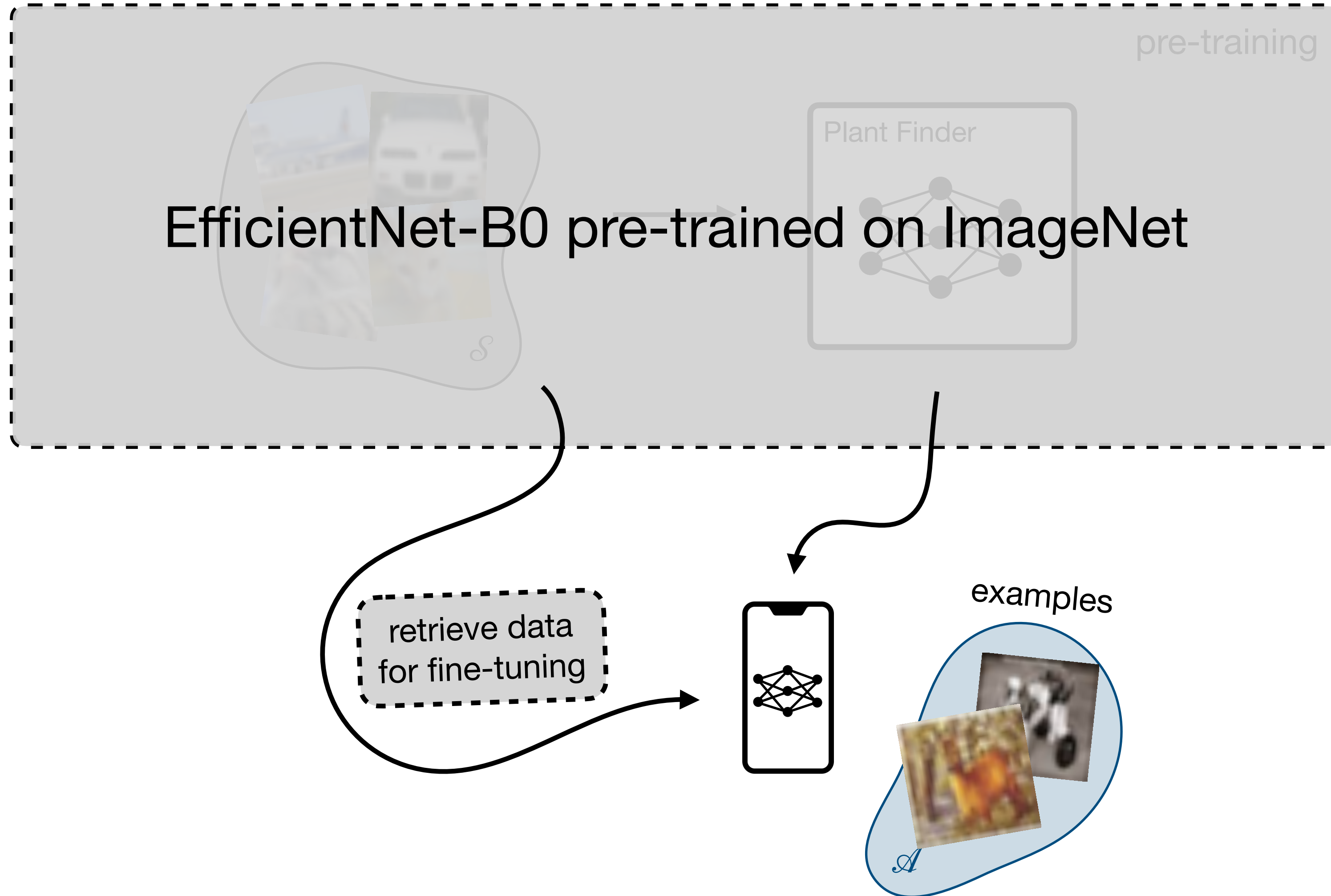
Fine-Tuning



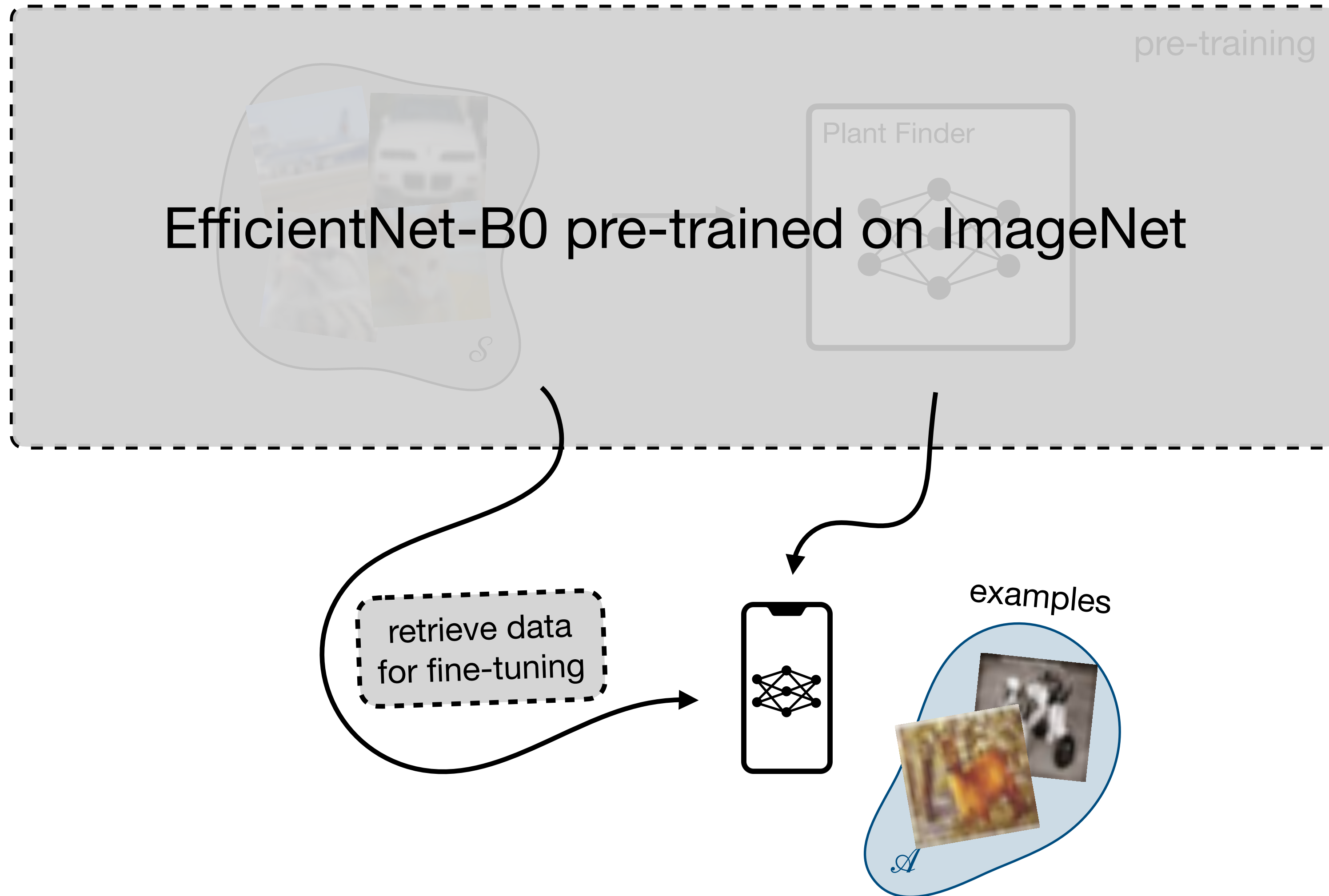
Fine-Tuning on CIFAR-100



Fine-Tuning on CIFAR-100

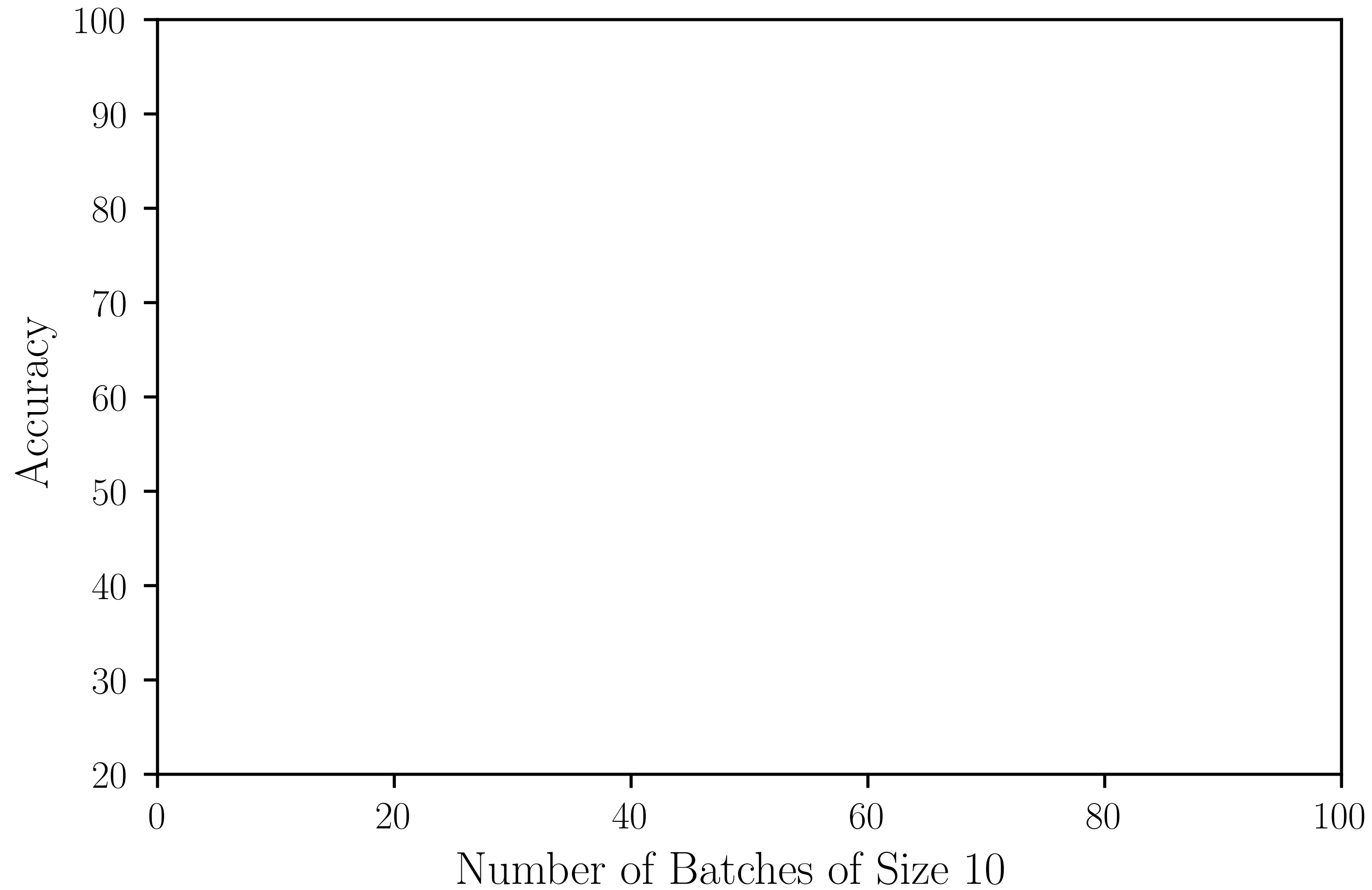


Fine-Tuning on CIFAR-100

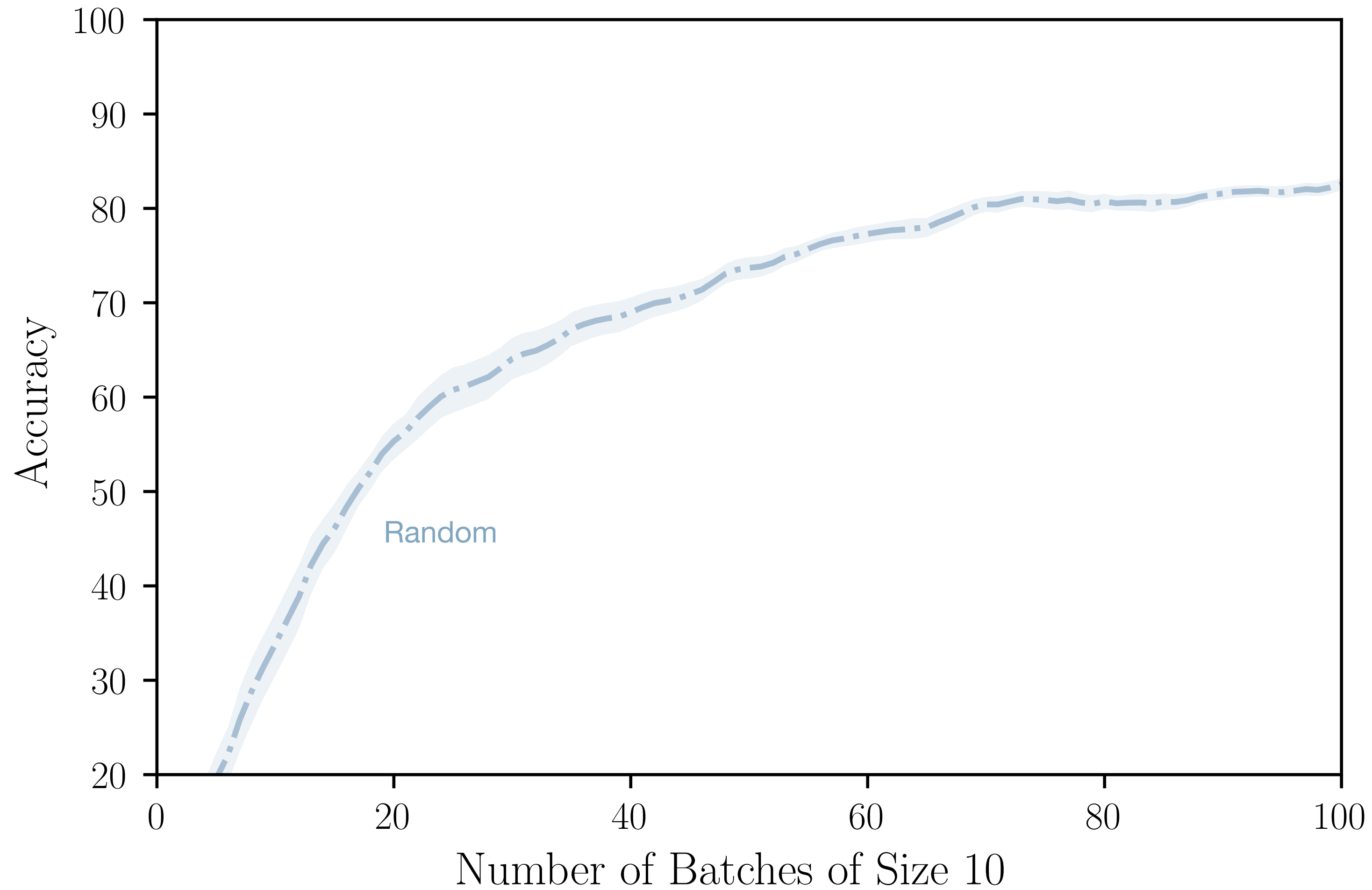


Goal: high accuracy on fresh examples from \mathcal{A}

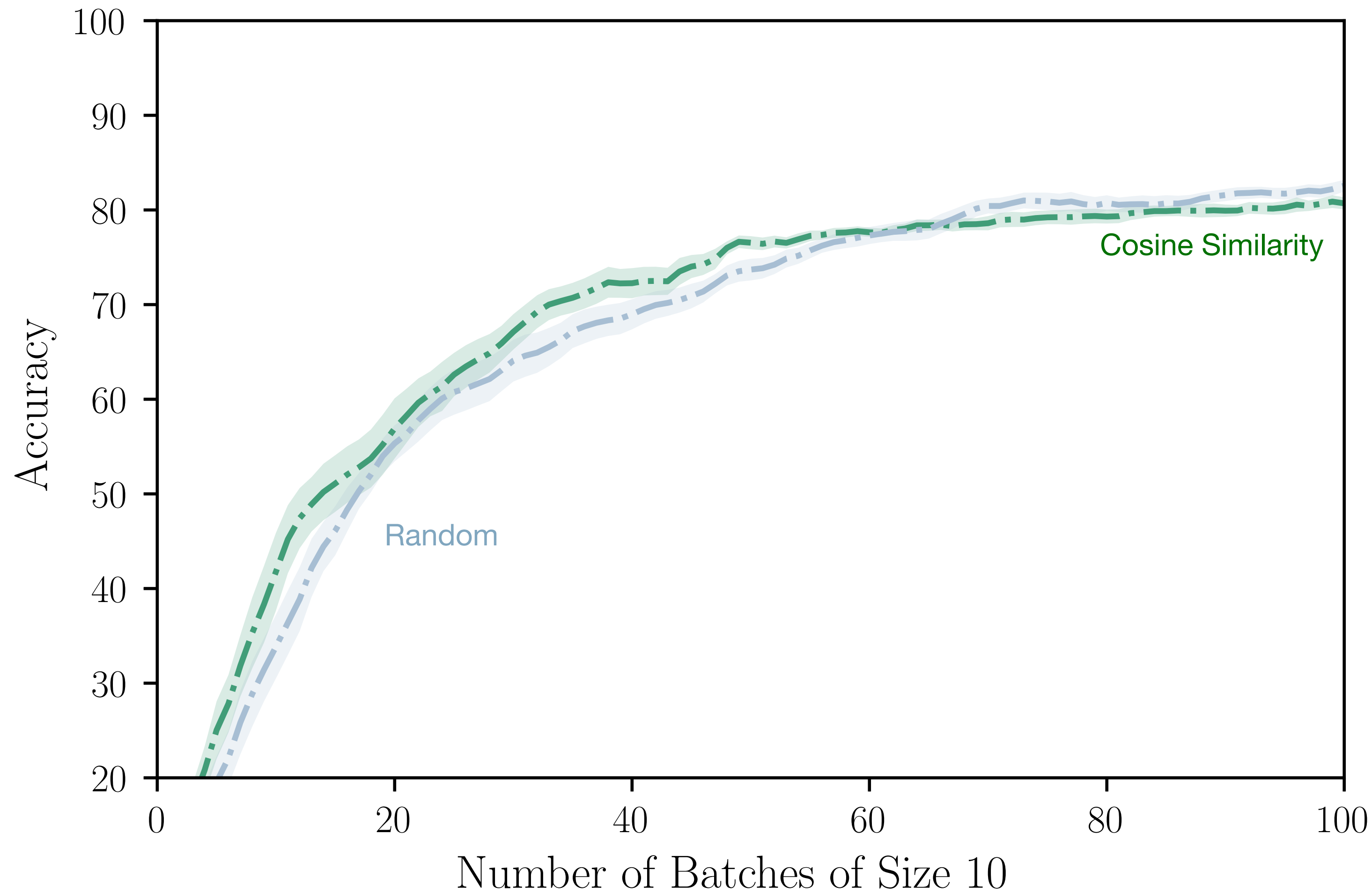
Fine-Tuning on CIFAR-100



Fine-Tuning on CIFAR-100



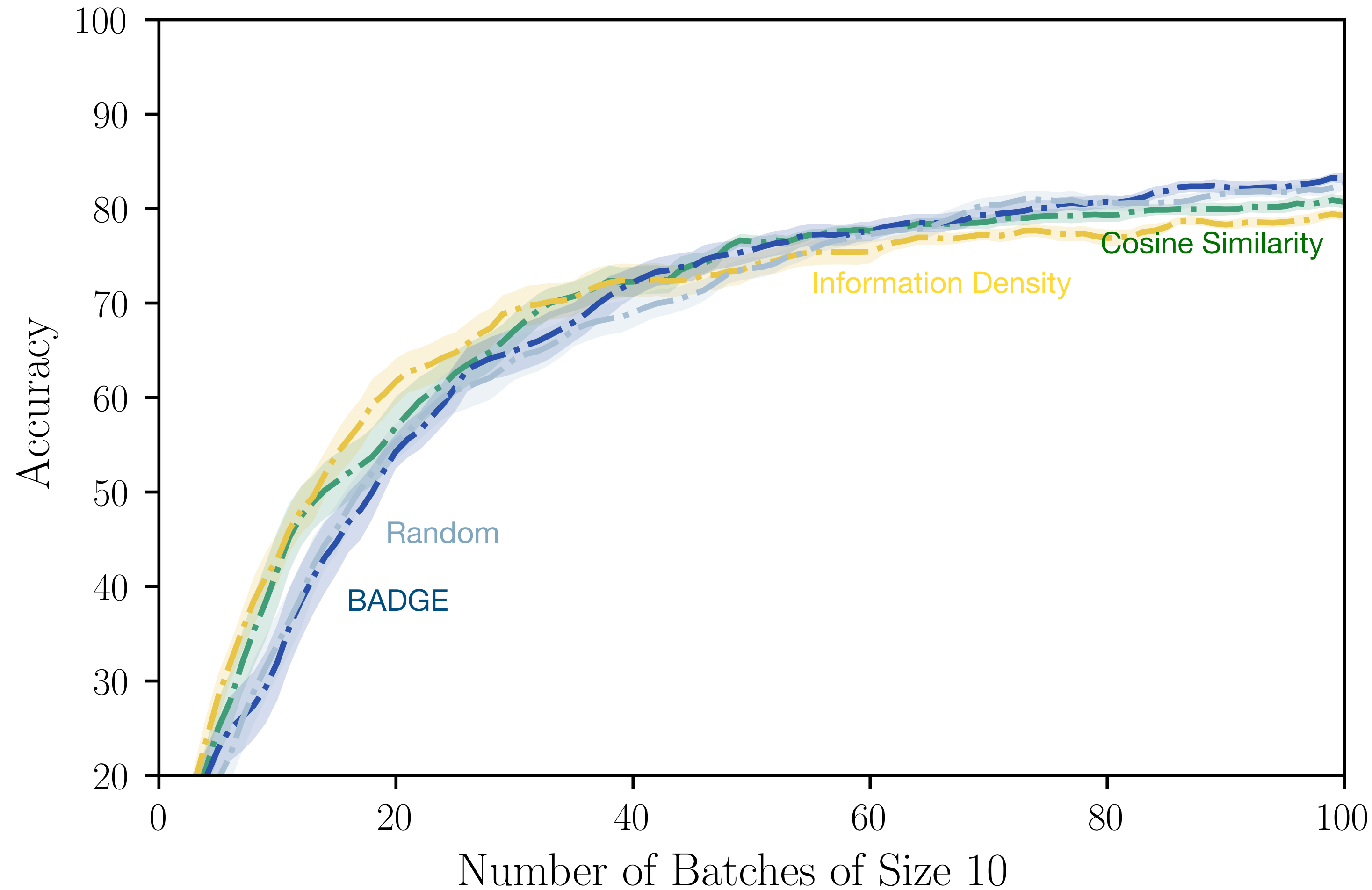
Fine-Tuning on CIFAR-100



Cosine Similarity: only relevance

$$\arg \max_{x \in \mathcal{S}} \frac{1}{|\mathcal{A}|} \sum_{x' \in \mathcal{A}} \underbrace{\Delta(\phi(x), \phi(x'))}_{\text{Cor}[f(x), f(x') | D_{n-1}]}$$

Fine-Tuning on CIFAR-100



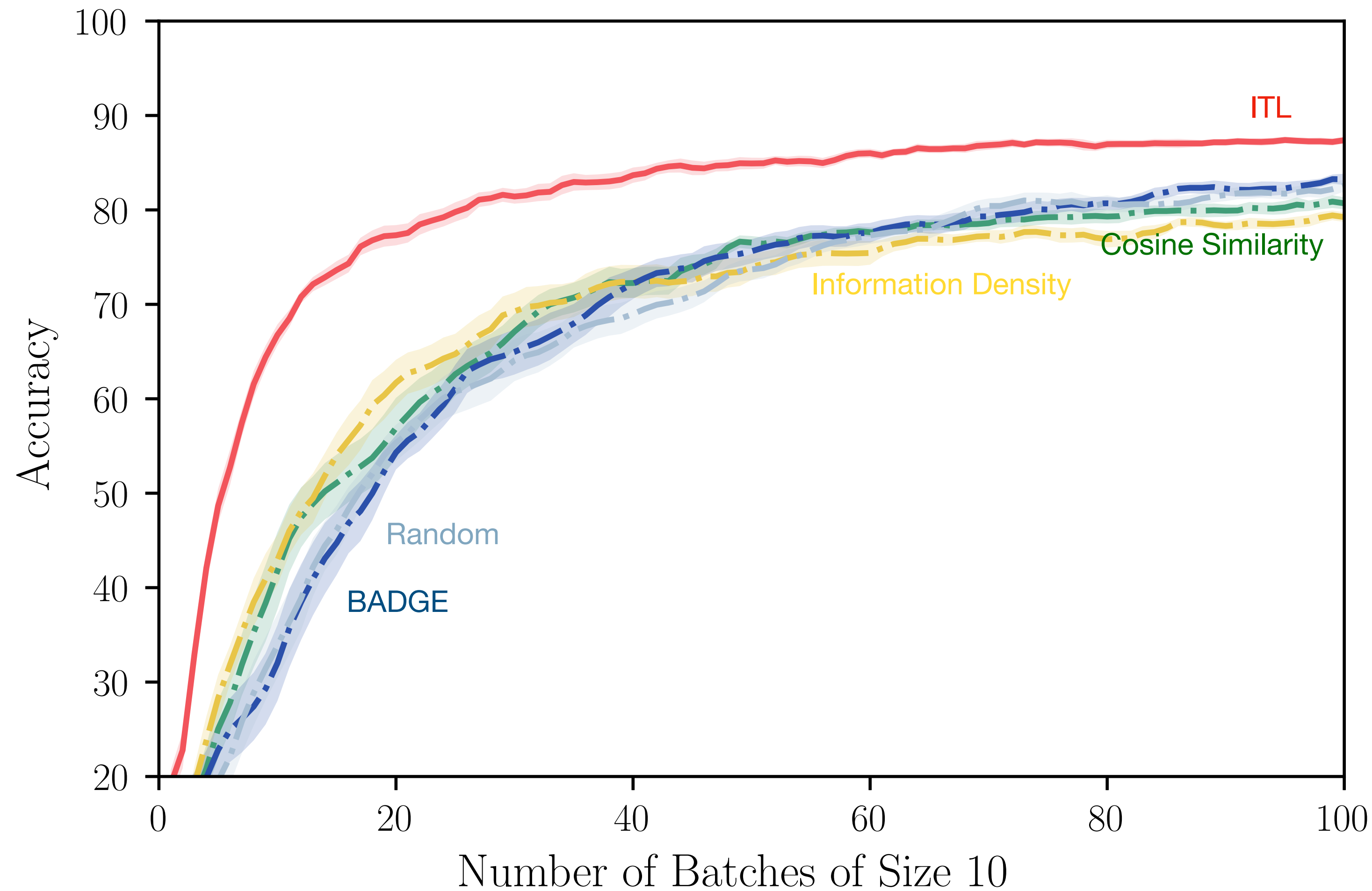
Cosine Similarity: only relevance

$$\arg \max_{x \in \mathcal{S}} \frac{1}{|\mathcal{A}|} \sum_{x' \in \mathcal{A}} \underbrace{\Delta(\phi(x), \phi(x'))}_{\text{Cor}[f(x), f(x') | D_{n-1}]}$$

Information Density:
only relevance

BADGE:
only diversity

Fine-Tuning on CIFAR-100



Cosine Similarity: only relevance

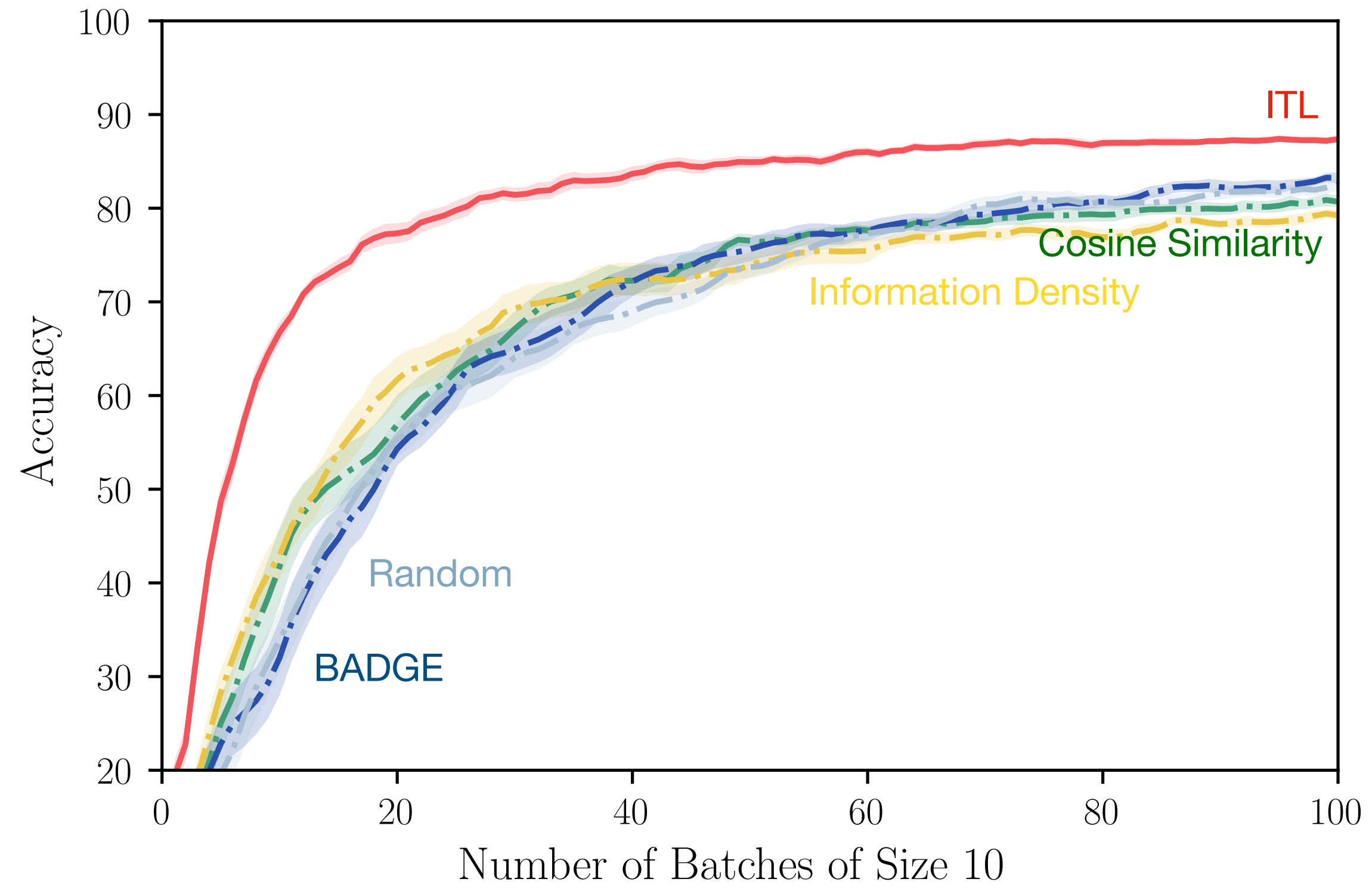
$$\arg \max_{x \in \mathcal{S}} \frac{1}{|\mathcal{A}|} \sum_{x' \in \mathcal{A}} \underbrace{\Delta(\phi(x), \phi(x'))}_{\text{Cor}[f(x), f(x') | D_{n-1}]}$$

Information Density:
only relevance

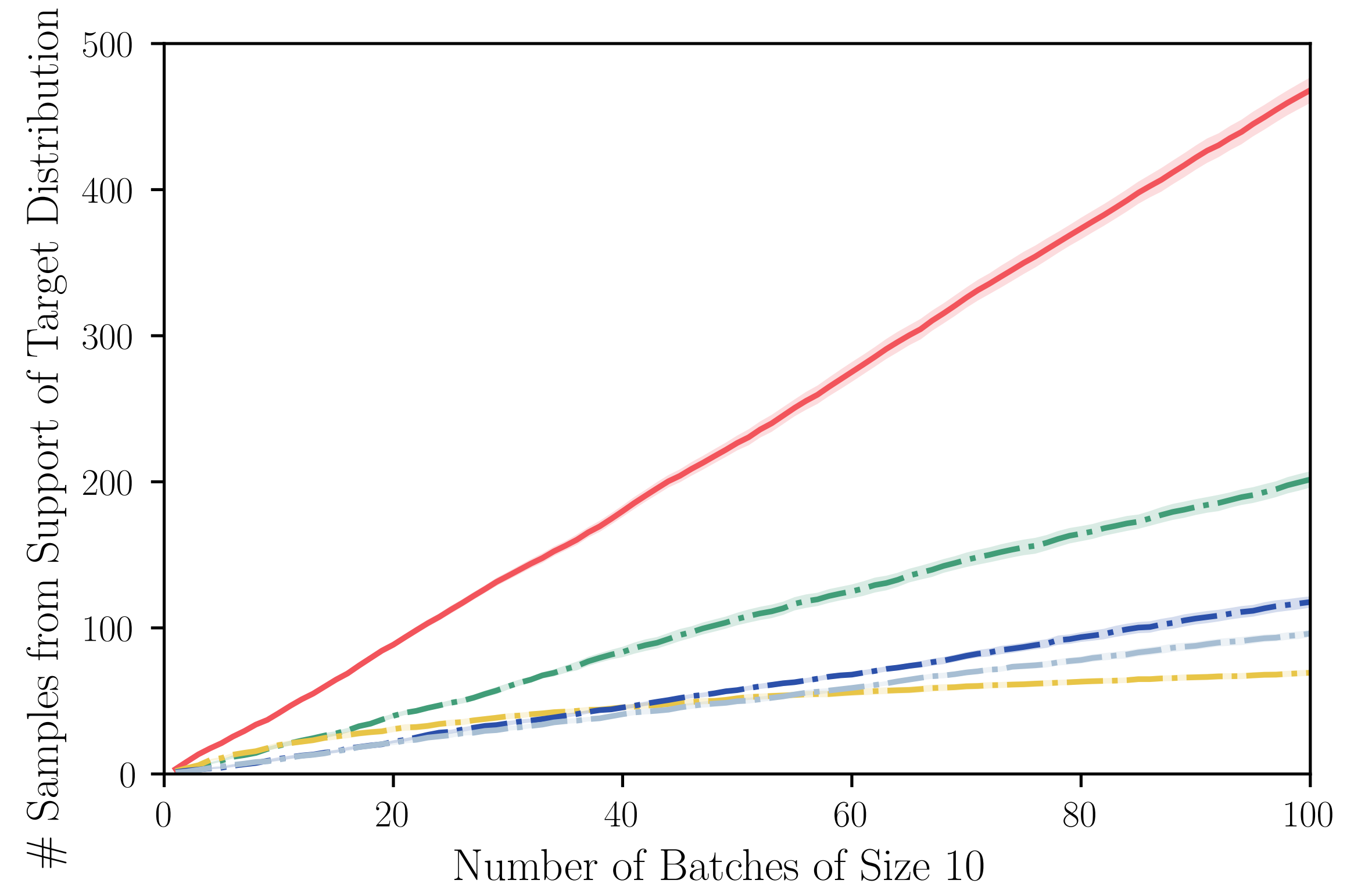
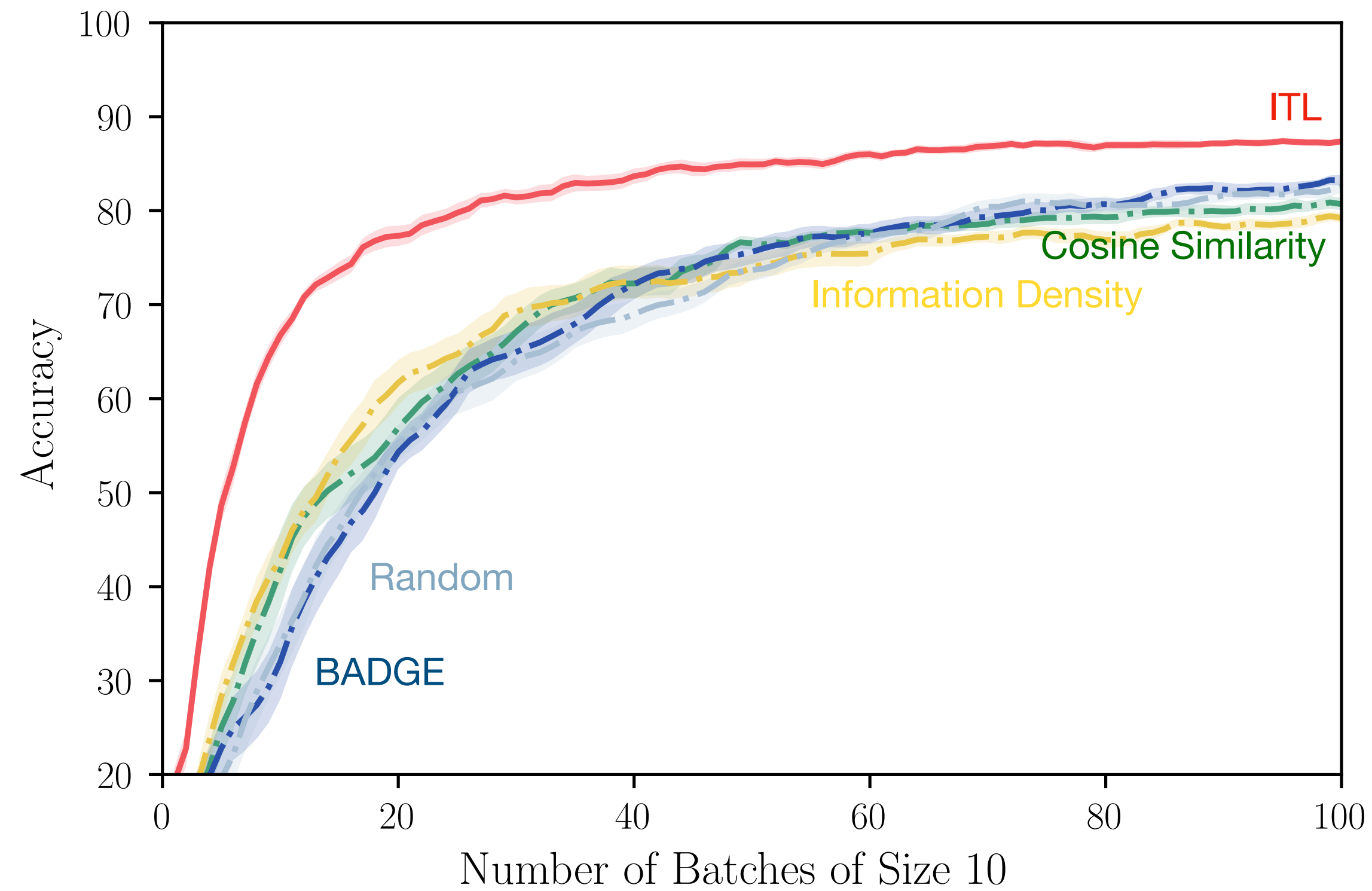
BADGE:
only diversity

ITL:
relevance + diversity

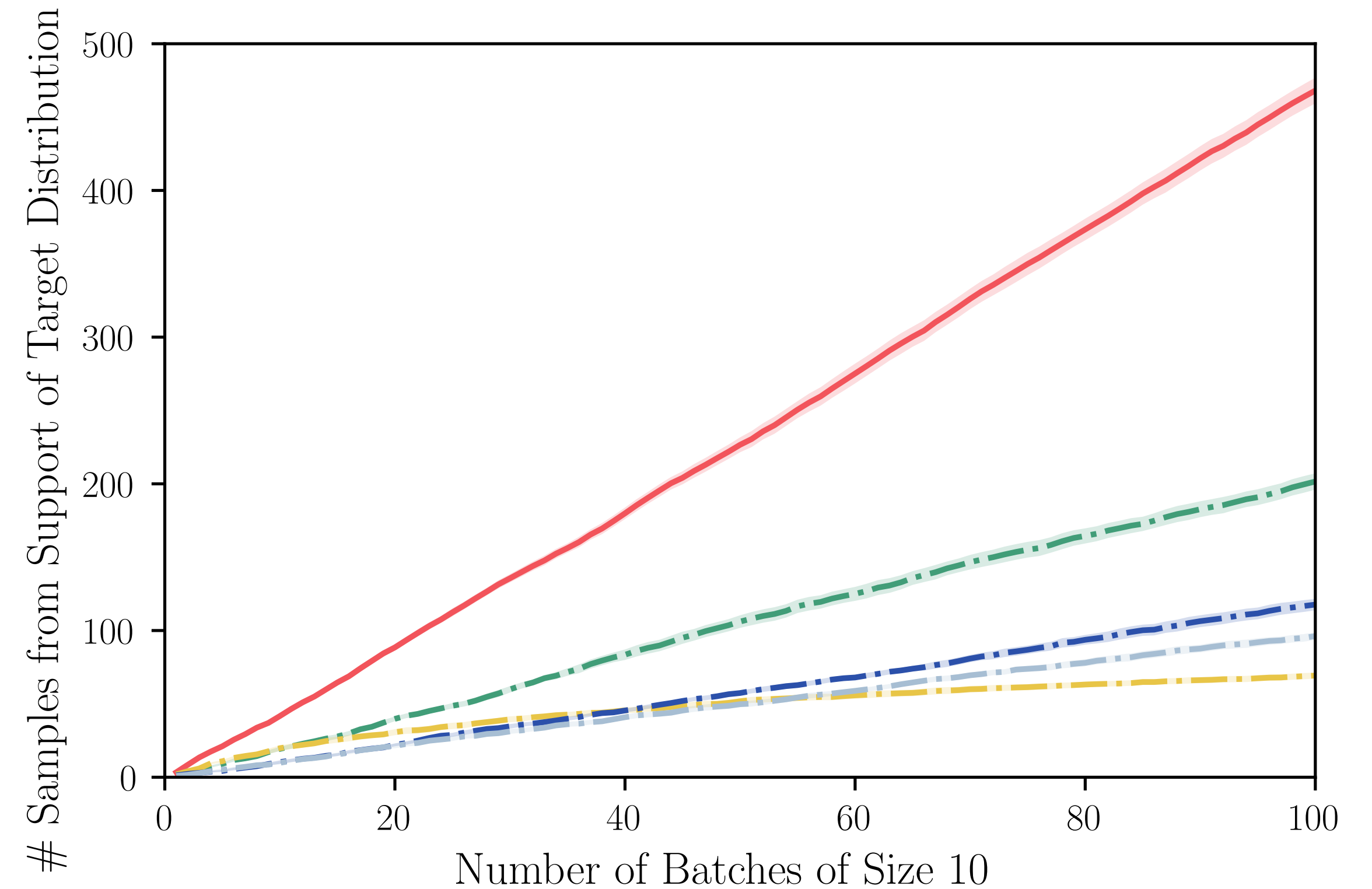
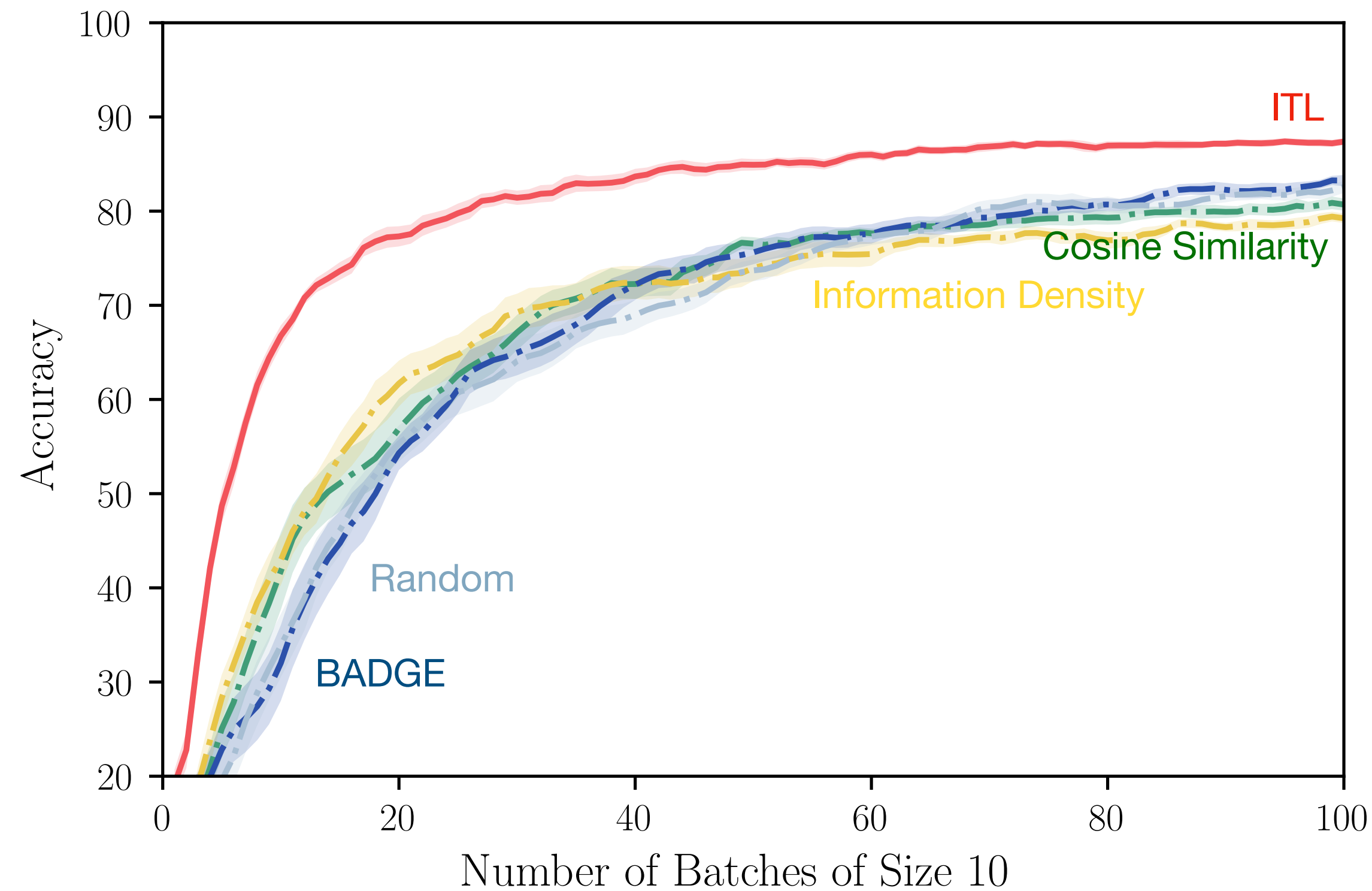
Fine-Tuning on CIFAR-100



Fine-Tuning on CIFAR-100



Fine-Tuning on CIFAR-100



ITL generalizes Cosine Similarity to query & batch sizes larger than 1!

Outlook



```
from afsl import ActiveDataLoader

train_loader = ActiveDataLoader.initialize(dataset, target, batch_size=32)

while not converged:
    batch = dataset[train_loader.next(model)]
    model.step(batch)
```

Outlook

```
from afsl import ActiveDataLoader

train_loader = ActiveDataLoader.initialize(dataset, target, batch_size=32)

while not converged:
    batch = dataset[train_loader.next(model)]
    model.step(batch)
```

- Fine-tuning in domains other than image classification on standard datasets
- Connection between learning and retrieval (in-context learning)
- Analyzing submodularity of retrieval / ITL
- Other applications of Transductive Active Learning (Safe BO, ...happy to chat!)

Bibliography

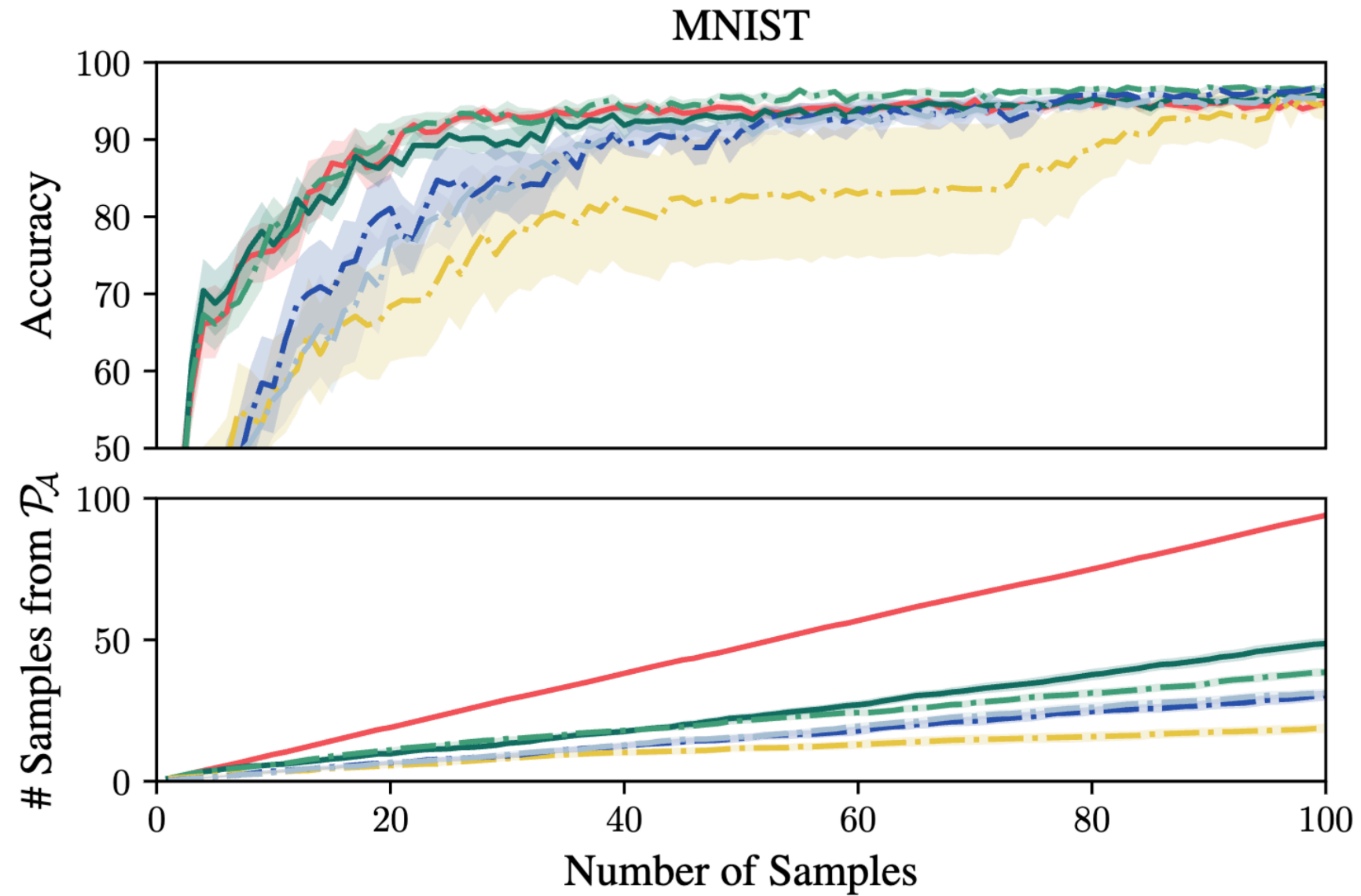
1. MacKay, D. J. Information-based objective functions for active data selection. *Neural computation*, 4(4), 1992.
2. Settles, B. and Craven, M. An analysis of active learning strategies for sequence labeling tasks. In *EMNLP*, 2008.
3. Ash, J. T., Zhang, C., Krishnamurthy, A., Langford, J., and Agarwal, A. Deep batch active learning by diverse, uncertain gradient lower bounds. *ICLR*, 2020.

Appendix

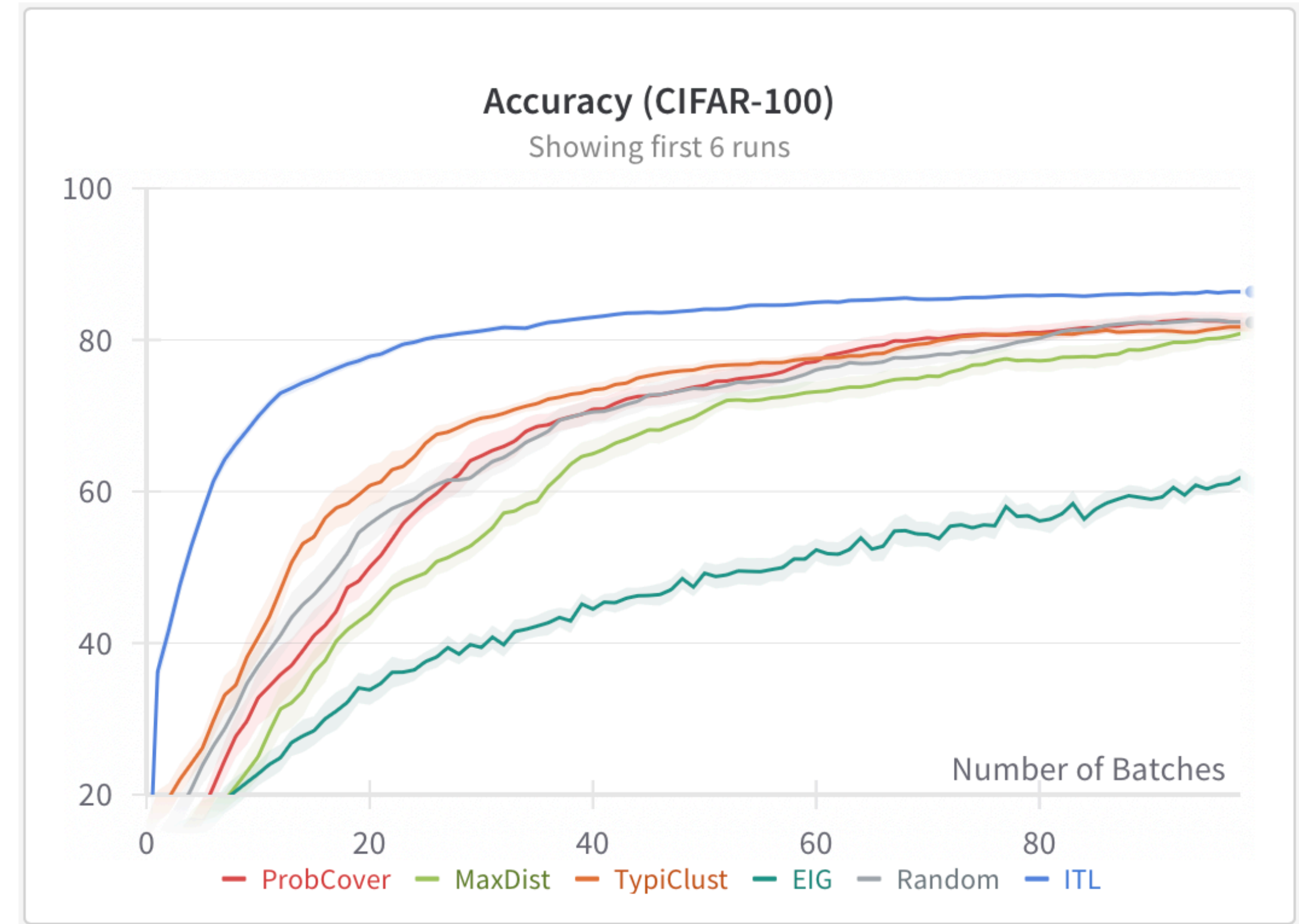
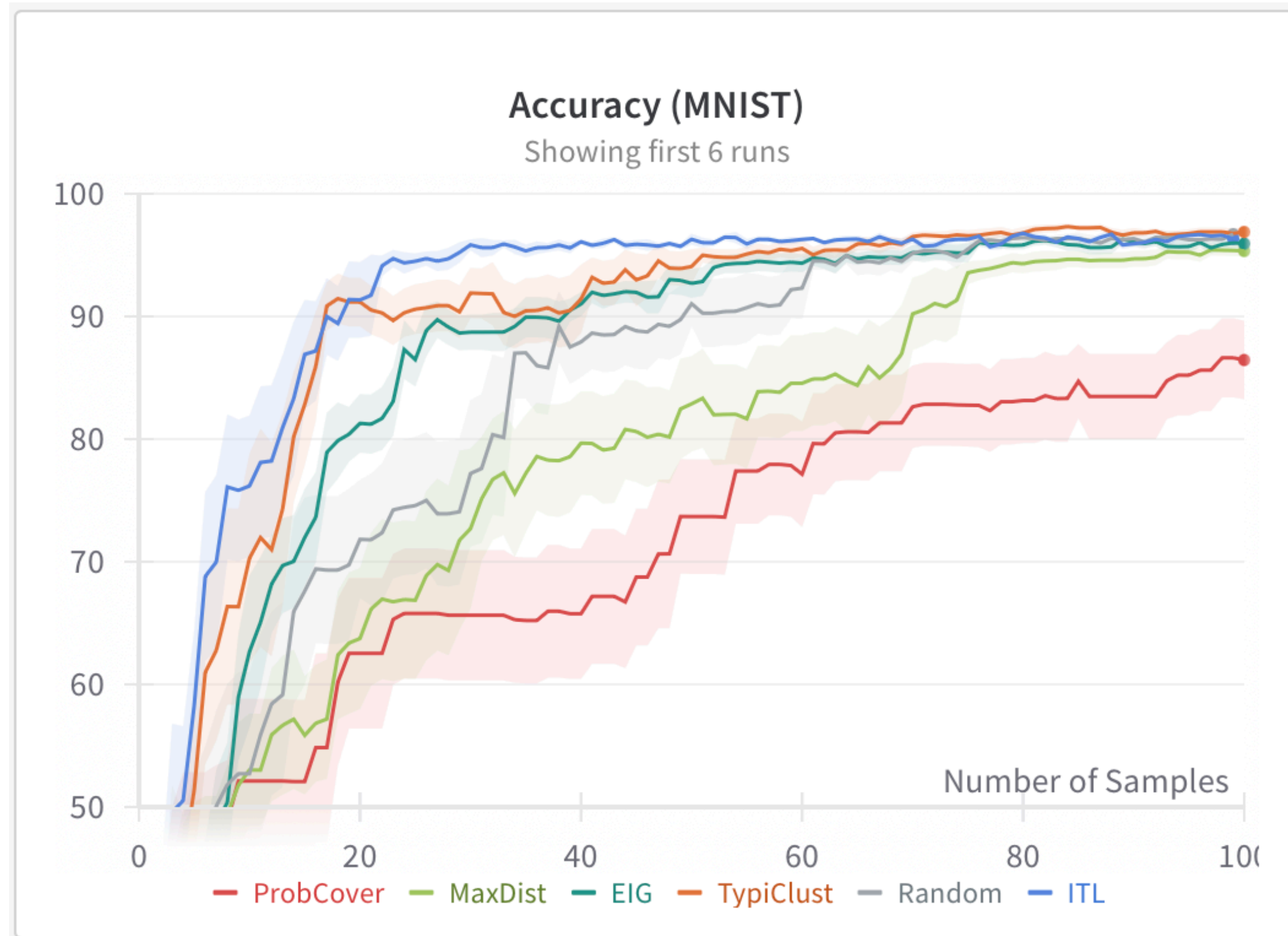
Embeddings

$$\phi(x) = \nabla_{\theta} \ell(f(x; \theta), \hat{y}(x)) \Big|_{\theta = \hat{\theta}}$$

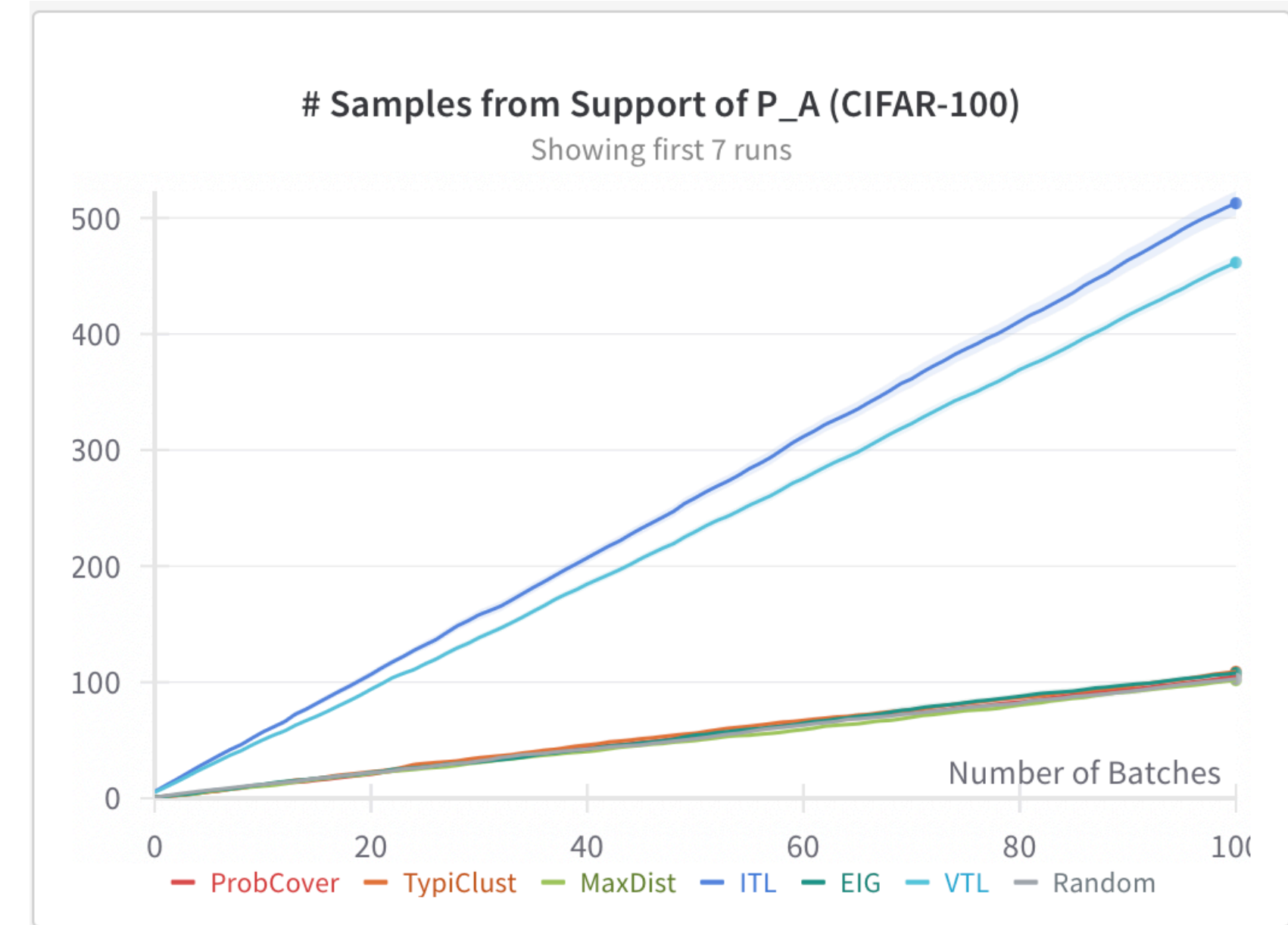
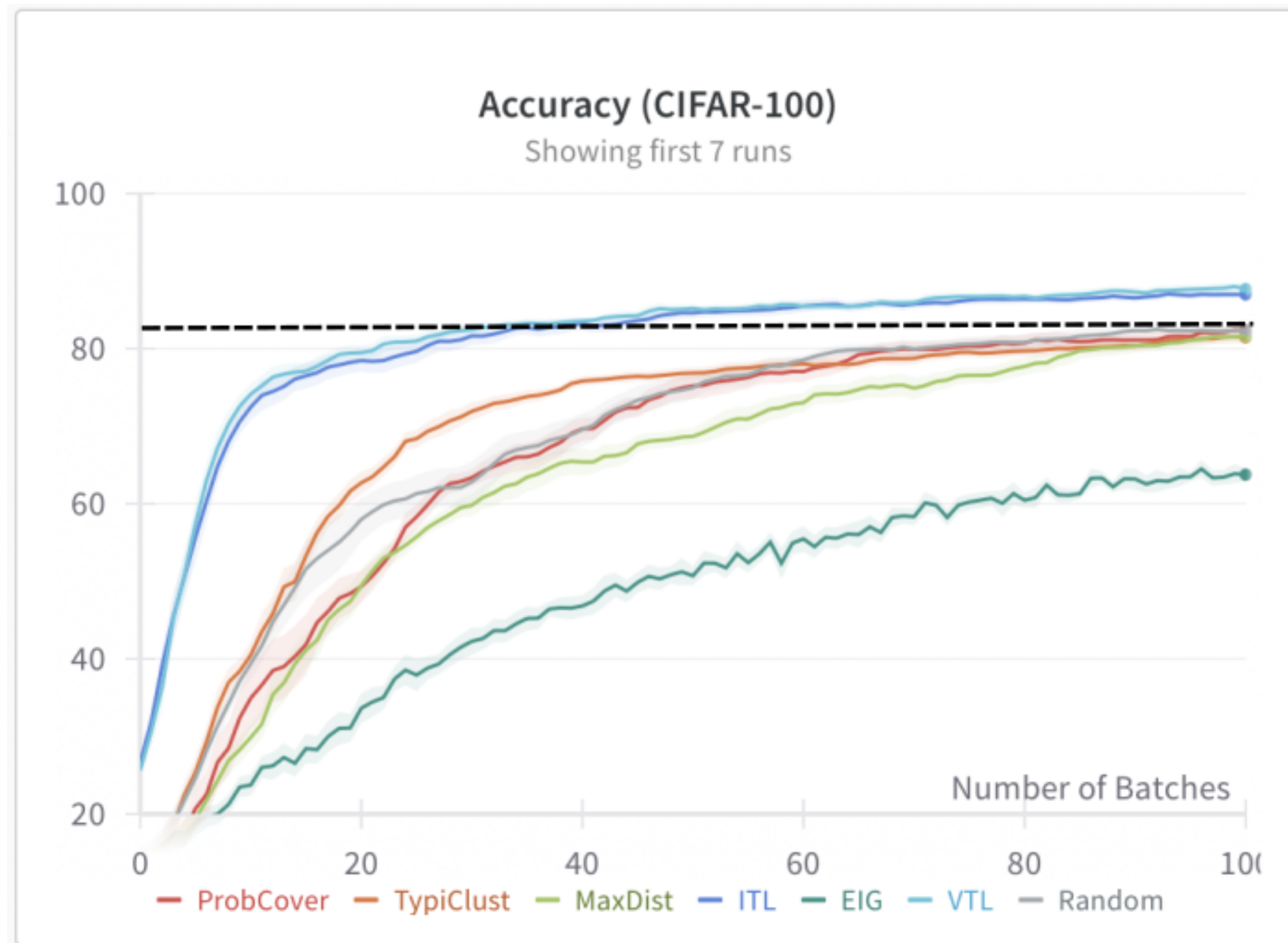
Without Pre-Training



Additional Baselines



Allow Sampling from Target Set: $\mathcal{A} \subseteq \mathcal{S}$



Setting where the sample space is $\mathcal{S} \cup \mathcal{A}$, i.e., includes the target space. The dashed black line is the accuracy after training on \mathcal{A} only where $|\mathcal{A}| = 100$.

Batch Selection via Conditional Embeddings

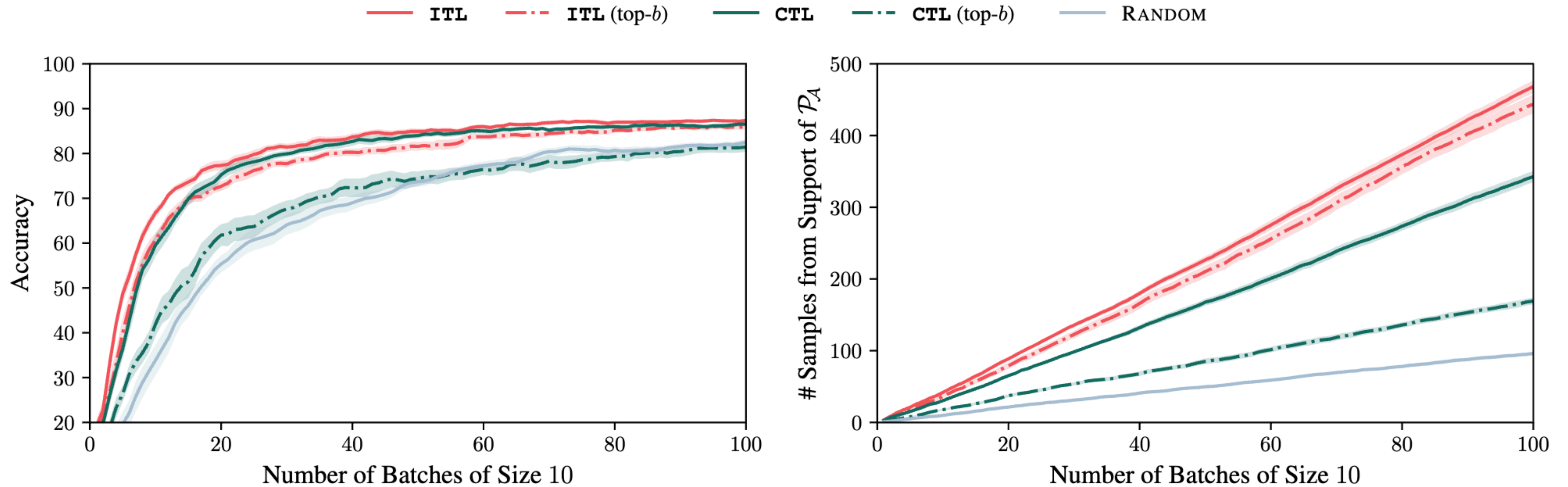


Figure 8. Advantage of batch selection via conditional embeddings over top- b selection in the CIFAR-100 experiment.