Task Discovery: Finding the Tasks that Neural Networks Generalize on

Andrei Atanov, Andrei Filatov, Teresa Yeo, Ajay Sohmshetty, Amir Zamir

taskdiscovery.epfl.ch



Task Discovery Problem

 Neural Networks can fit any task, i.e., any labeling of a set of images [1,2], but what tasks can they generalize on?

Task Discovery: finds different generalizable tasks automatically.

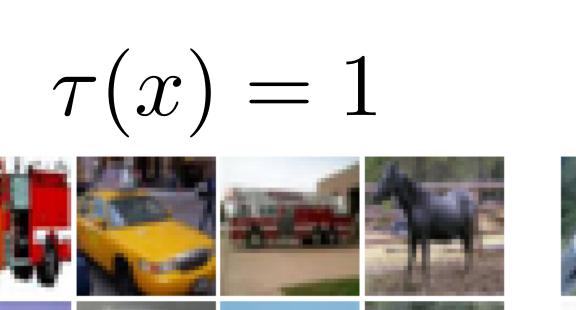
- How do these tasks look? What do they indicate?
- Generalizable tasks reflect the inductive biases of NNs ⇒ can help us <u>understand deep learning</u> better.

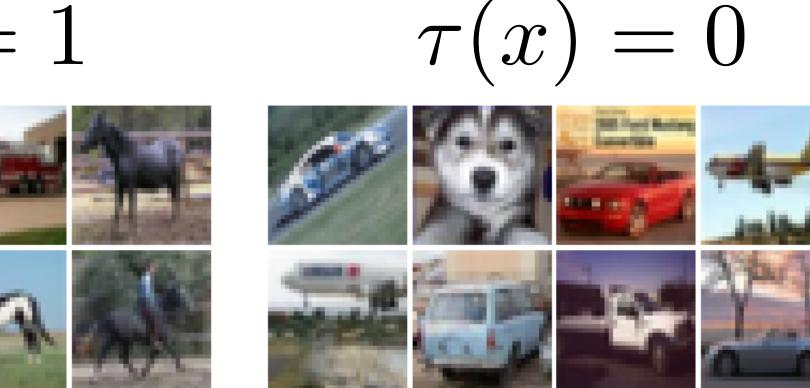
Generalizability via <u>Agreement Score</u>

 How do we differentiate between generalizable (e.g., human-labeled) and non-generalizable (e.g., random-labeled) tasks <u>quantifiably</u>?



A Human-Labeled Task





labels are assigned randomly to each image

A Random-Labeled Task

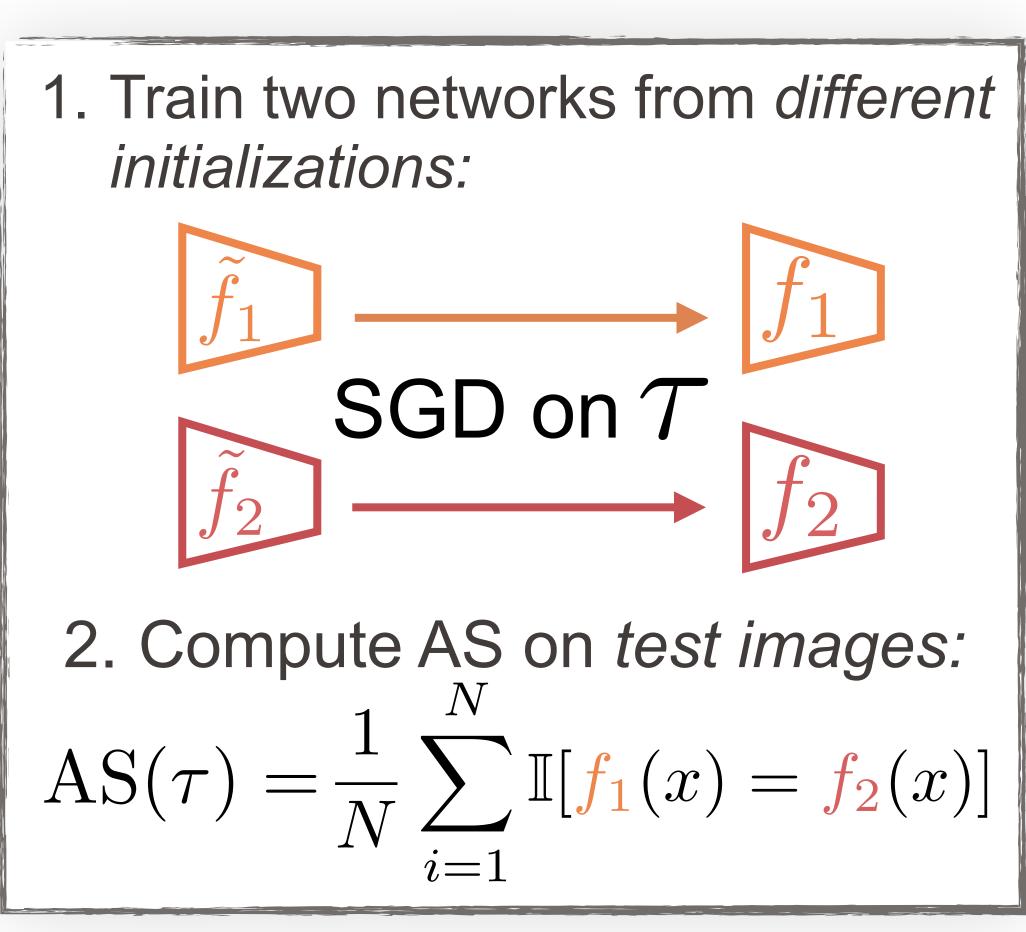
Tasks:

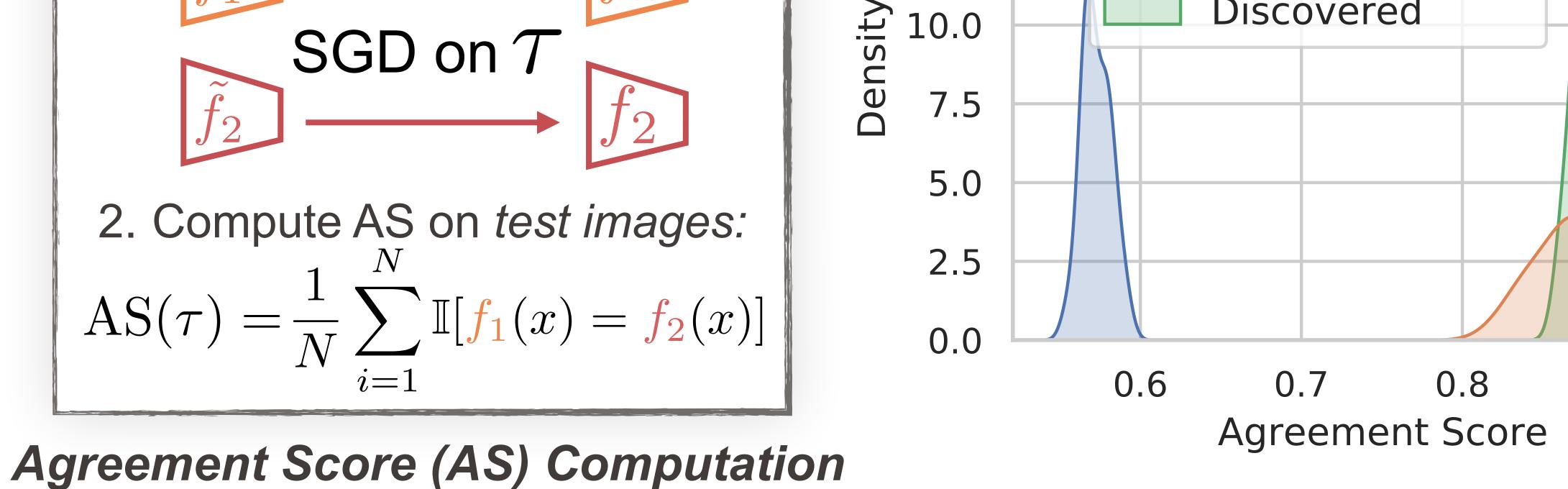
Random-labelled

Discovered

Human-labelled

 We use Agreement Score (AS), a generalization-based quantity, to identify generalizable tasks computationally:



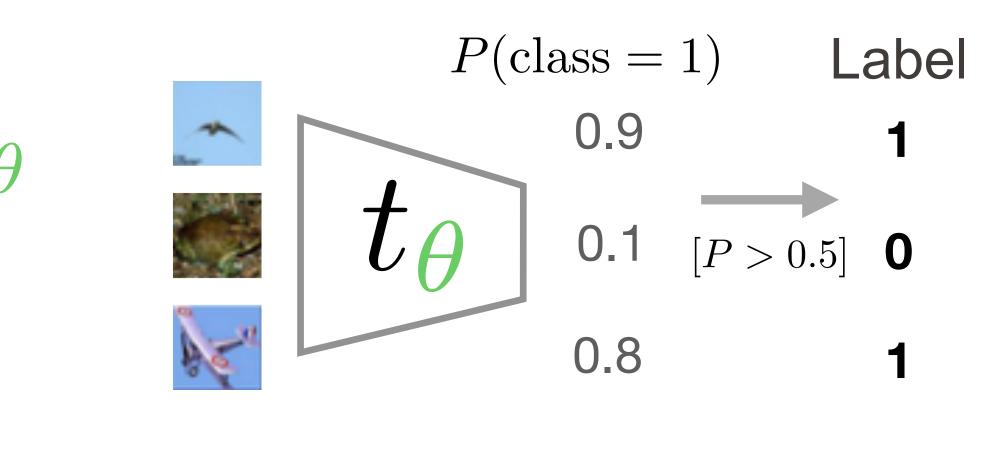


Automatic Discovery via Meta-Optimization

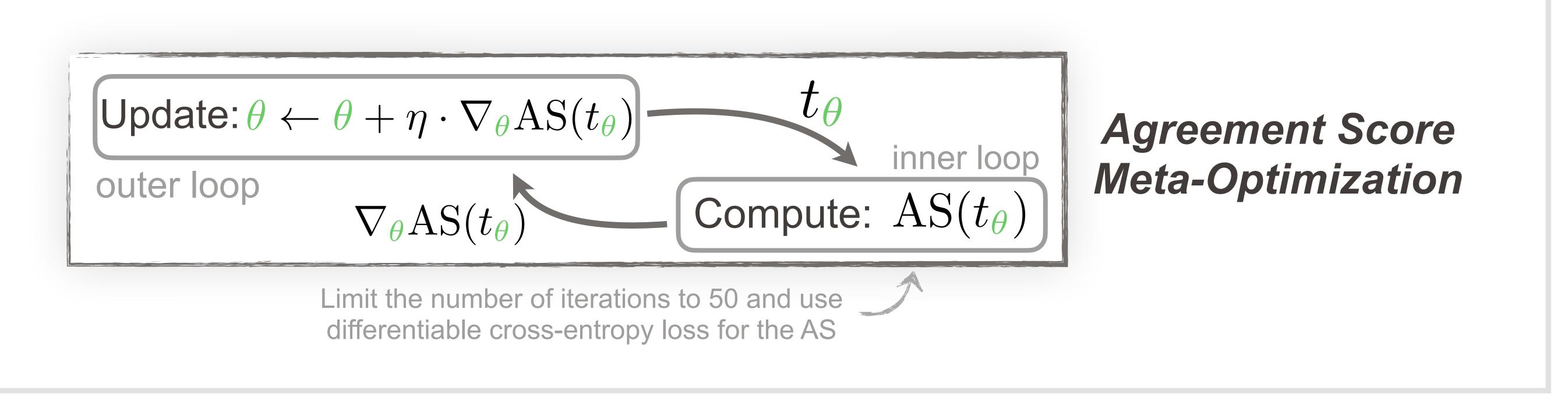
 We optimize the labeling of a set of images, resulting in a high-AS discovered task:

$$\tau_{\mathrm{d}} \leftarrow \operatorname*{arg\,max} \mathsf{AS}(\tau)$$

- Model the labeling with a <u>task network</u> $t_{ heta}$ that outputs the label for each image.

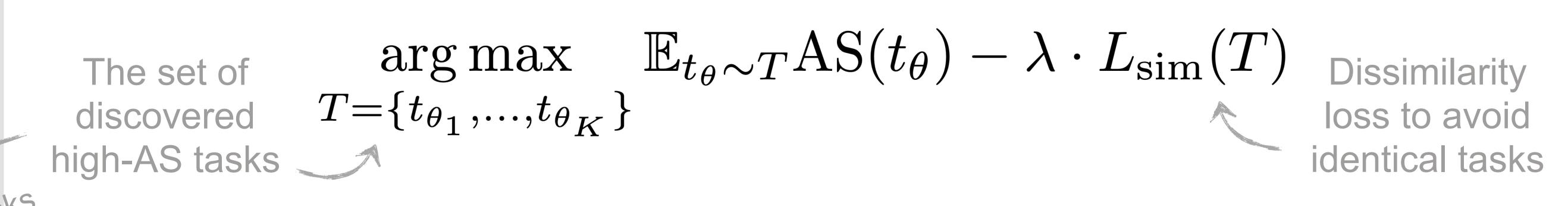


- Use meta-optimization:
- Outer loop: update the task network that provides labels
- Inner loop: compute the AS of a given task

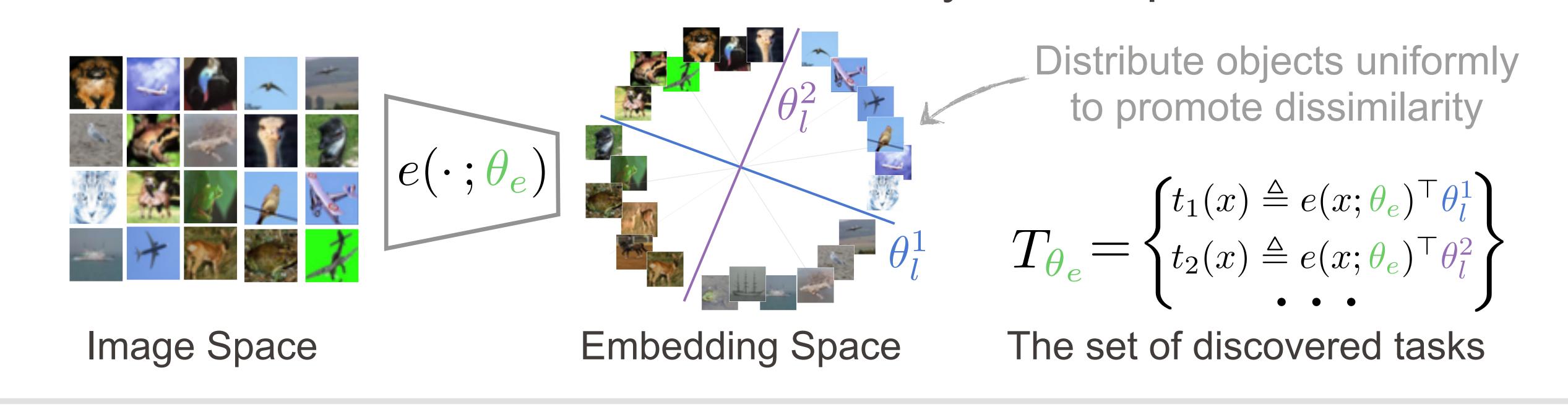


Find <u>Different Tasks</u> via Dissimilarity Constraint

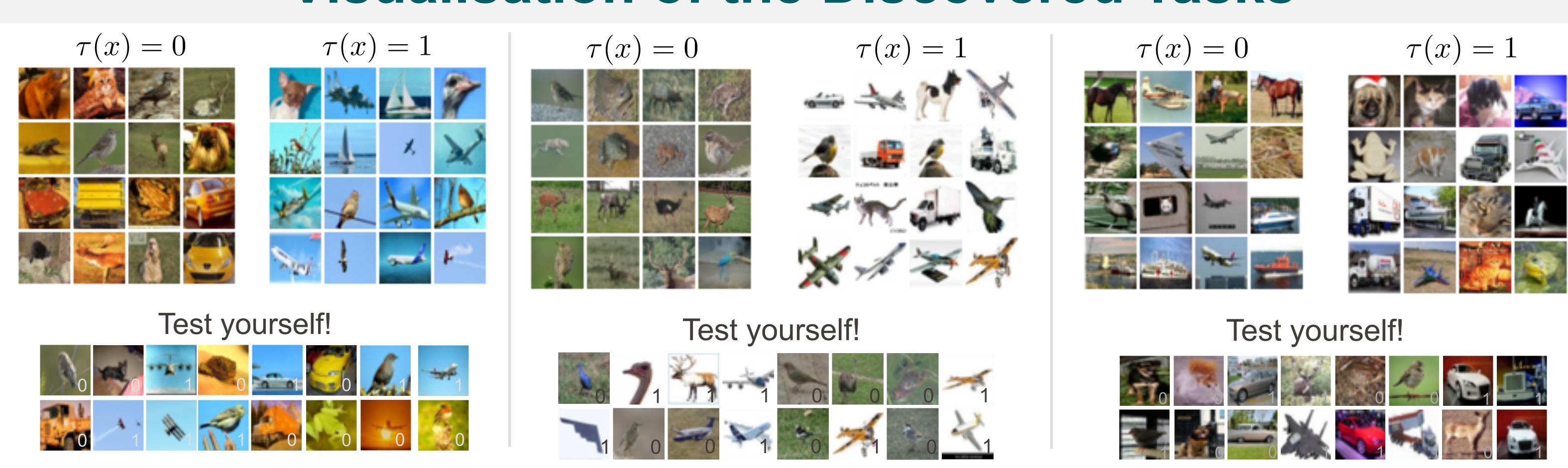
- Meta-optimization results in a single task, but there could be many high-AS tasks we want to discover.
- We aim to discover <u>a set</u> of dissimilar tasks:



Use a shared encoder with a uniformity loss in practice:



Visualisation of the Discovered Tasks

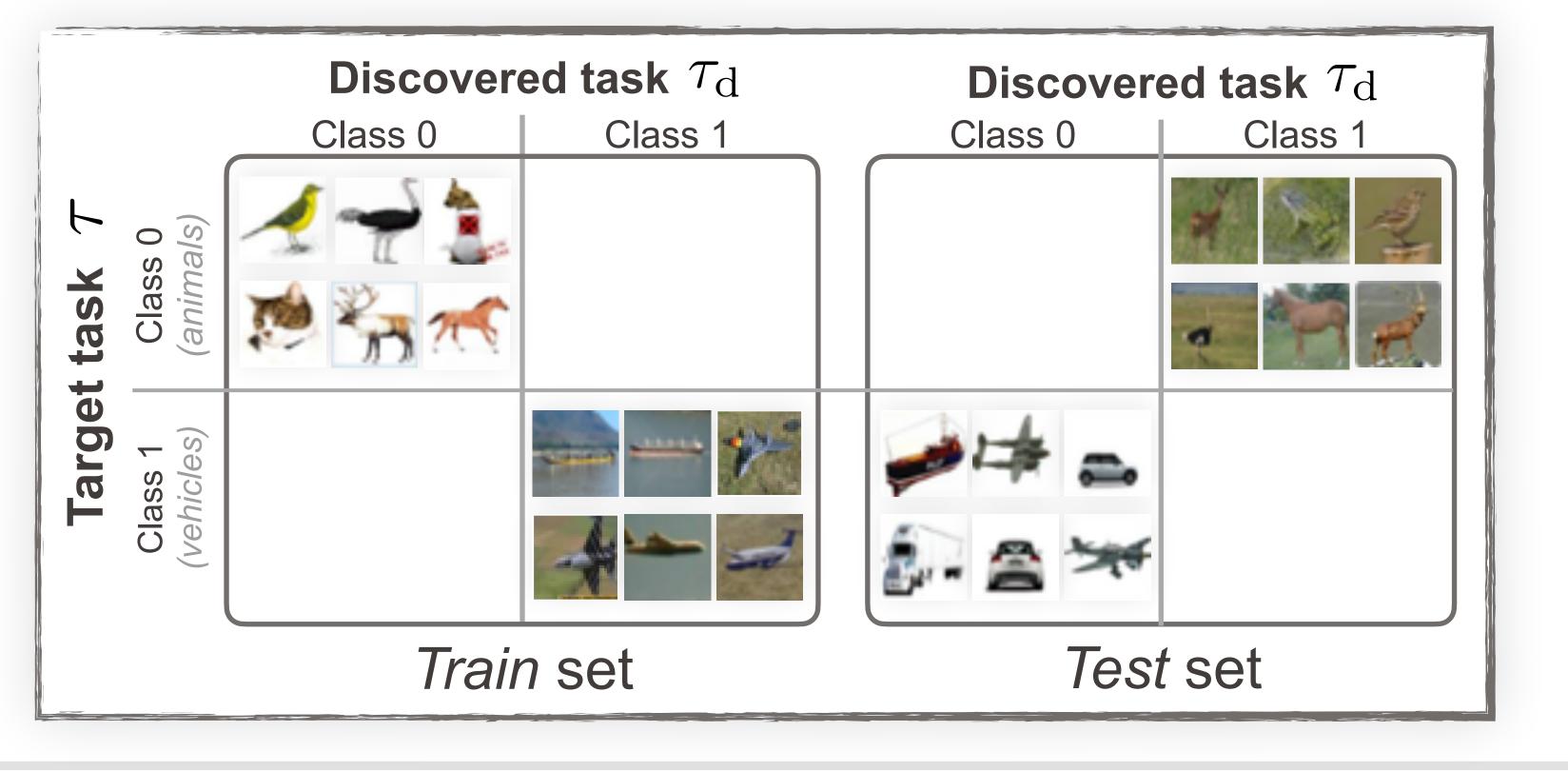


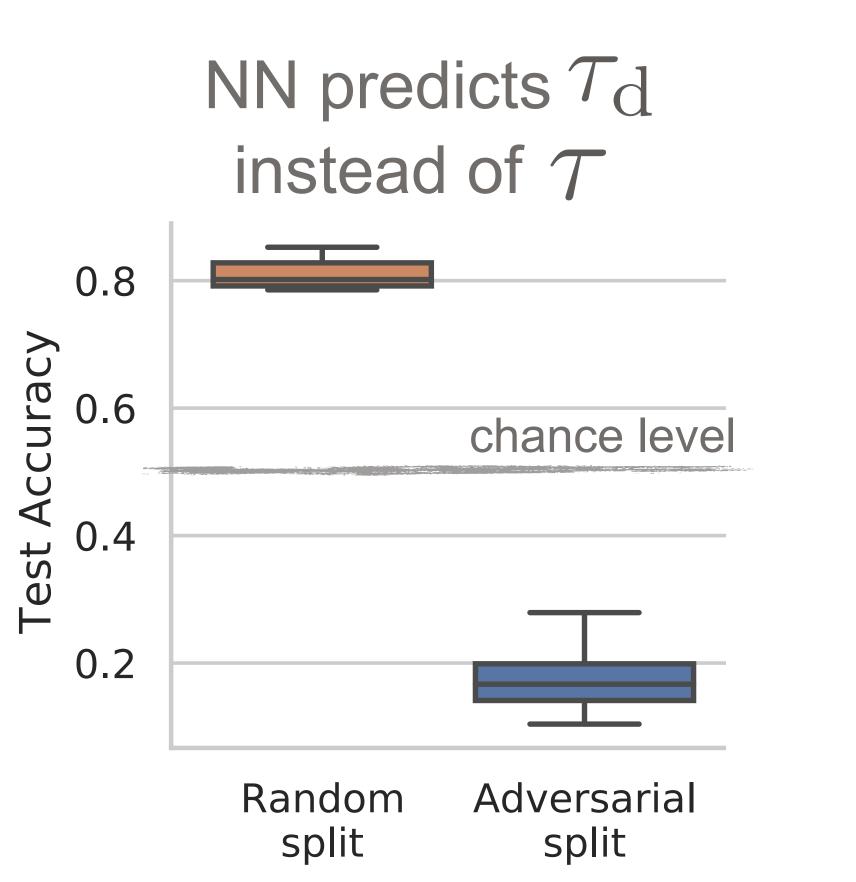
 Some tasks are not easy to interpret visually by humans, and are not meant to be! (see Sec. 5.4 of the paper)

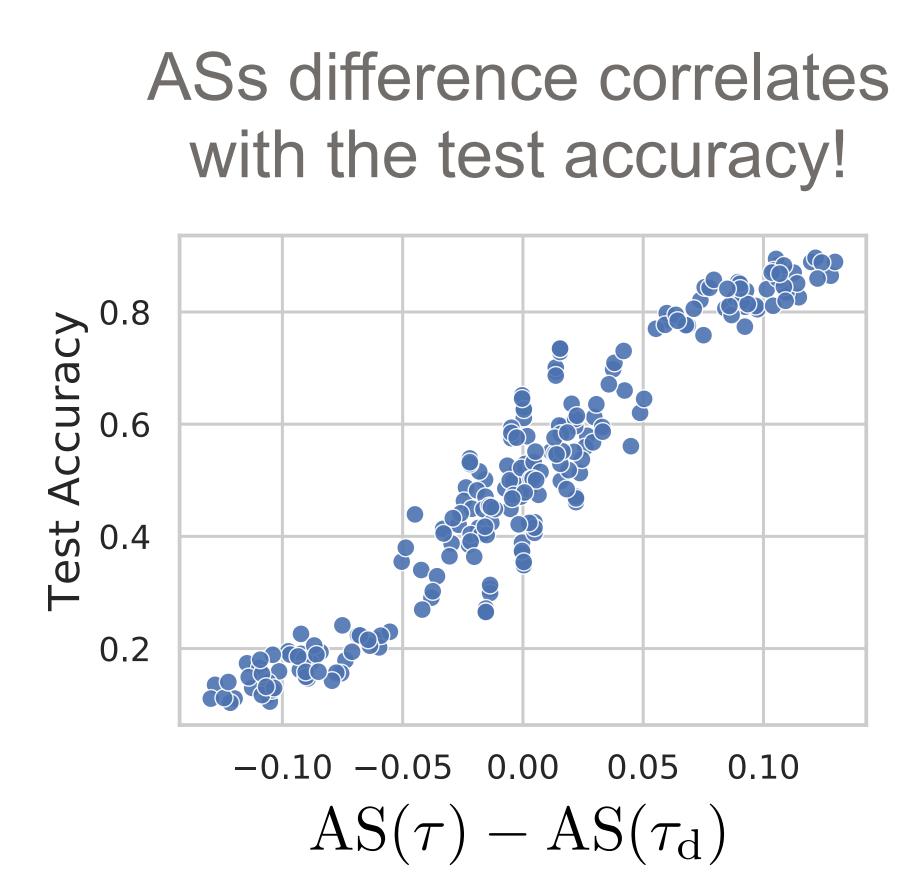
Discovered Tasks Reveal NNs failure modes via Adversarial Splits

- Discovered tasks reflect inductive biases of NNs and statistics patterns in data.
- Adversarial splits use them to <u>highlight NNs' failure modes</u> by creating a <u>spurious</u> correlation between a discovered and target to "fool" the network:

Adversarial train-test split







Summary

- Find generalizable tasks automatically via AS meta-optimization.
- Discovered tasks reflect NNs' inductive biases and statistical patterns in data \Rightarrow can help us analyze and understand them better.
 - Example: <u>adversarial splits</u> reveal NNs' failure modes.

References:

- [1] Understanding deep learning requires rethinking generalization, Zhang et al., ICLR 2017
- [2] Small ReLU networks are powerful memorizers: a tight analysis of memorization capacity, Yun et al., NeurlPS 2019