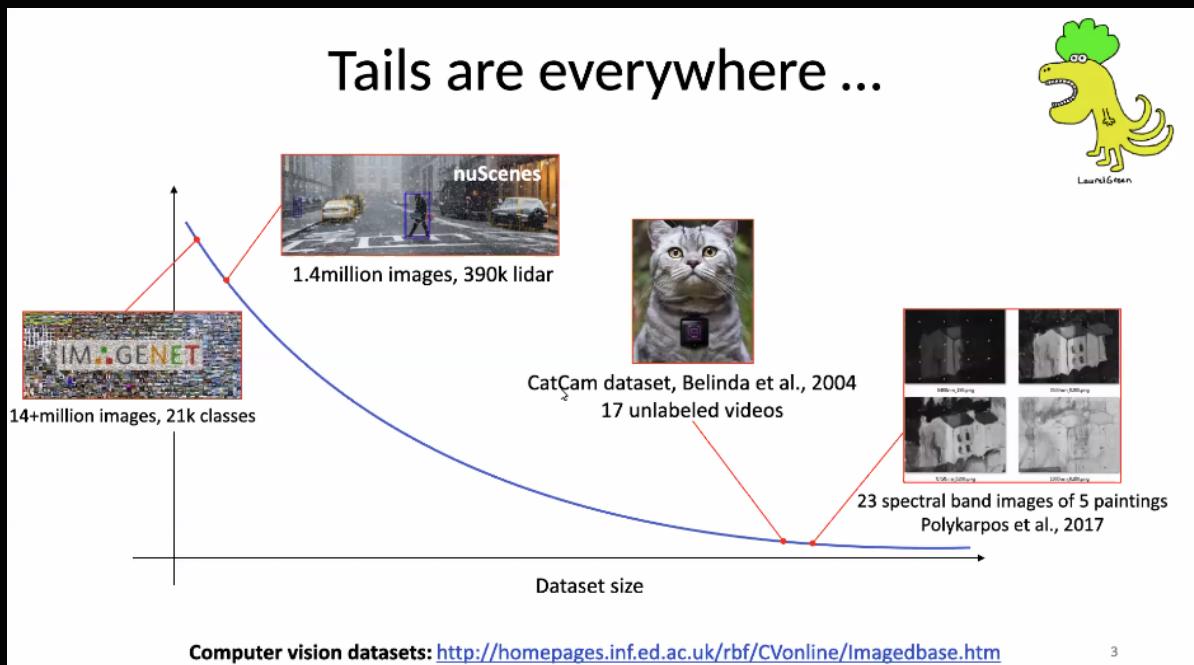
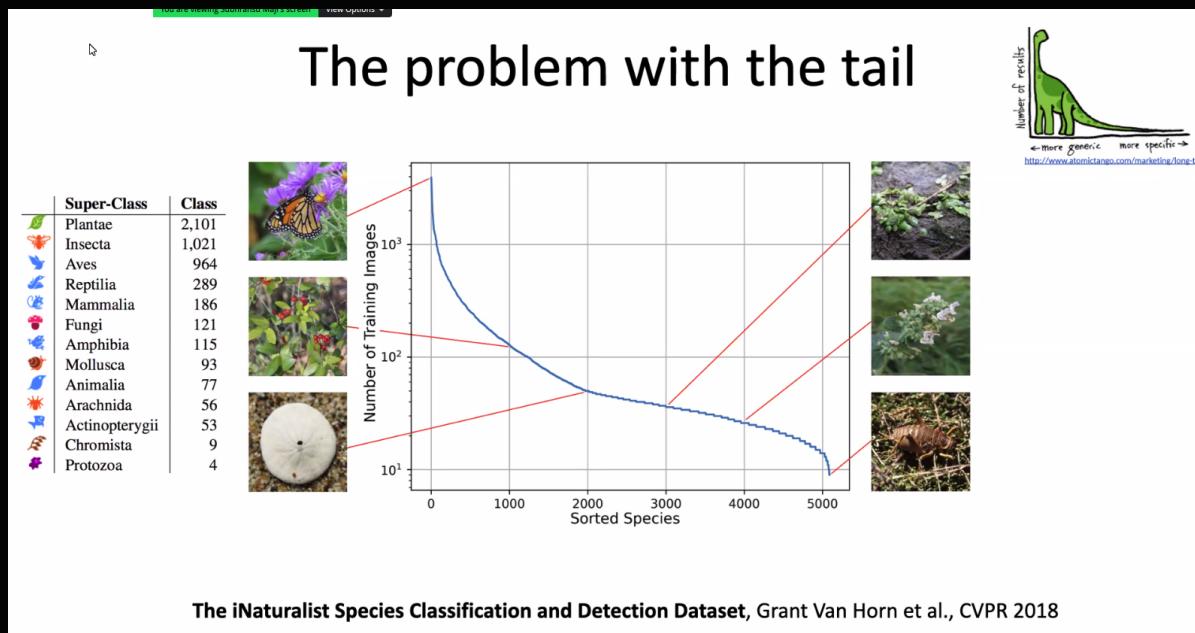


Day 15:

Speaker: Prof Subhranshu Maji , University of Massachusetts.

Title: Learning from Little data



Transfer learning to the rescue

Models trained on ImageNet transfer well to other vision tasks

- What about other pre-trained models (e.g., iNaturalist)?
- What about novel domains (e.g., medical images)?

We lack a general framework to reason about tasks and domains

Why?

- Few-shot learning (e.g., recommend the right expert)
- Multi-tasking (e.g., which tasks are mutually beneficial)
- Semi/self-supervised learning on which dataset?

4

Talk outline

I. **Task2Vec:** Vector representations of tasks [ICCV'19]

- Captures task similarity, task difficulty, etc.
- Useful for model recommendation



Alessandro Achille

II. When does self-supervision improve few-shot learning? [CVPR'20]

III. A realistic-evaluation of semi-supervised learning [CVPR'21]

- Self/semi—supervision is effective as long as the unlabeled images are from a similar domain
- Requires a notion of domain similarity



Jong-Chyi Su Zezhou Cheng

Task Embedding

Alessandro, Michael, Rahul, Avinash, Subhransu, Charless, Stefano, Pietro @ ICCV'19

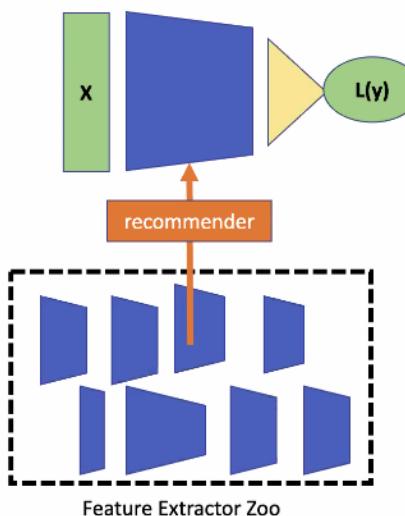


- What are similar tasks?
- What architecture should I use?
- What pre-training dataset?
- What hyper parameters?
- Do I need more training data?
- How difficult is this task?

If we have a universal **vectorial representation** of **tasks** we can frame all sorts of interesting computer vision application engineering problems as meta-learning problems.

6

Model recommendation



Brute Force:

Input: Task = (dataset, loss)

For each feature extractor architecture F :

1. Train classifier on $F(\text{dataset})$
2. Compute validation performance

Output: best performing model

Task recommendation:

Input: Task = (dataset, loss)

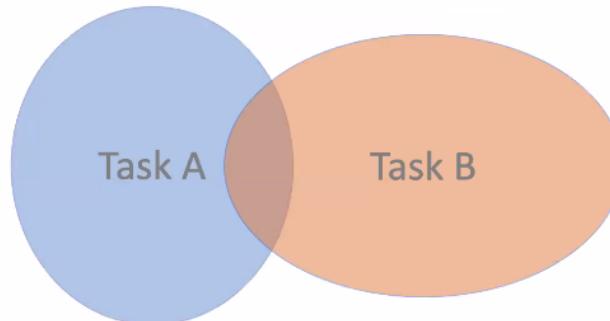
1. Compute task embedding $t = E(\text{Task})$
2. Predict best extractor $F = M(t)$
3. Train classifier on $F(\text{dataset})$
4. Compute validation performance

Output: best performing model

Train ahead of time.

7

Similarity measures on the space of tasks



Task = {images, labels, loss}

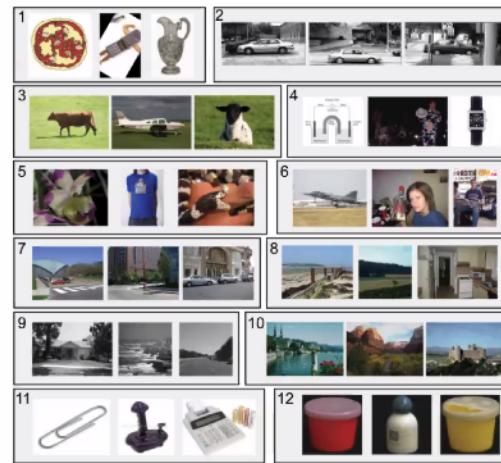
Dataset = {images, labels}

8

Similarity measures on the space of tasks

Domain similarity

Unbiased look at dataset bias, Torralba and Efros, CVPR 11



Caltech101 Tiny
MSRC Corel
UIUC PASCAL 07

LabelMe 15 Scenes
COIL-100 Caltech256
ImageNet SUN09

Similarity measures on the space of tasks

Domain similarity

Range / label similarity

- e.g., Taxonomic distance

$$D_{\text{tax}}(t_a, t_b) = \min_{i \in S_a, j \in S_b} d(i, j),$$

$D(\text{bird task}, \text{mammal task}) < D(\text{bird task}, \text{worm task})$

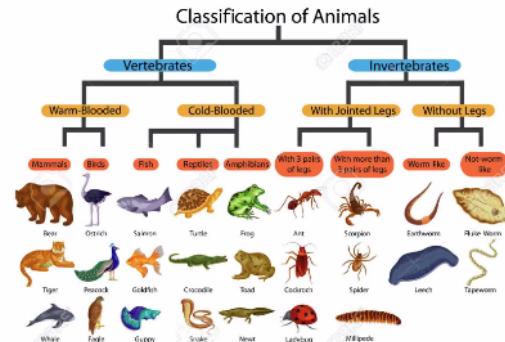


image source: https://www.123rf.com/photo_80712766_stock-vector-education-chart-of-biology-for-classification-of-animals-diagram.html

10

Similarity measures on the space of tasks

Domain similarity

Range / label similarity

- e.g., Taxonomic distance

$$D_{\text{tax}}(t_a, t_b) = \min_{i \in S_a, j \in S_b} d(i, j),$$

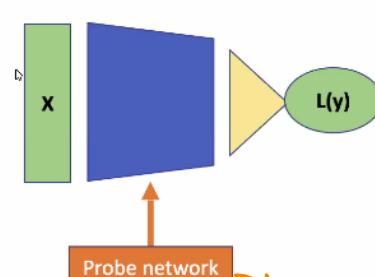
Transfer “distance”

- Fine-tune on task a followed by b

$$D_{\text{ft}}(t_a \rightarrow t_b) = \frac{\mathbb{E}[\ell_{a \rightarrow b}] - \mathbb{E}[\ell_b]}{\mathbb{E}[\ell_b]}$$

Task embedding using a probe network

- Given a **task**, train a classifier with the **task loss** on features from a generic “probe network”
- Compute gradients of **probe network** parameters w.r.t. task loss
- Use statistics of the probe parameter **gradients** as the fixed dimensional **task embedding**



say (any general net like AlexNet, ResNet etc)
find out which part of P.N is useful for solving the task

11

Task embedding using Fisher Information

- Given a **task**, train a classifier with the **task loss** on features from a generic “probe network”
- Compute gradients of **probe network** parameters w.r.t. task loss
- Use statistics of the probe parameter **gradients** as the fixed dimensional **task embedding**

$$F = \frac{1}{N} \sum_{i=1}^N \nabla \log p(x_i|\theta) \nabla \log p(x_i|\theta)^T$$

Intuition: F provides information about the **sensitivity** of the task performance to small perturbations of **parameters** in the probe network

$$\mathbb{E}_{x \sim \hat{p}} KL p_{\theta'}(y|x)p_{\theta}(y|x) = \delta\theta \cdot F \cdot \delta\theta + o(\delta\theta^2),$$

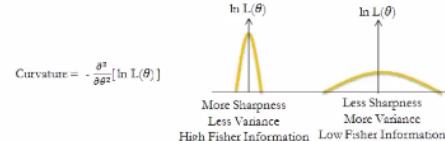


Figure source: <https://www.gaussianwaves.com/2012/10/score-and-fisher-information-estimator-sensitivity/>

13

Robust Fisher Computation

- For realistic CV tasks we want to use deep CNNs (e.g., ResNet) and estimate FIM for all the parameters.
- Challenge:** FIM can be hard to estimate (noisy loss landscape; high dimensions; small training set)
- Robust FIM**
 - Restrict it to a diagonal
 - Restrict it a single value per filter (CNN layer)
 - Robust estimation via perturbation

Estimate Λ of a Gaussian perturbation:

$$L(\hat{w}; \Lambda) = \mathbb{E}_{w \sim \mathcal{N}(\hat{w}, \Lambda)} [H_{p_w, \hat{p}} p(y|x)] + \beta KL(\mathcal{N}(0, \Lambda) \| \mathcal{N}(0, \lambda^2 I))$$

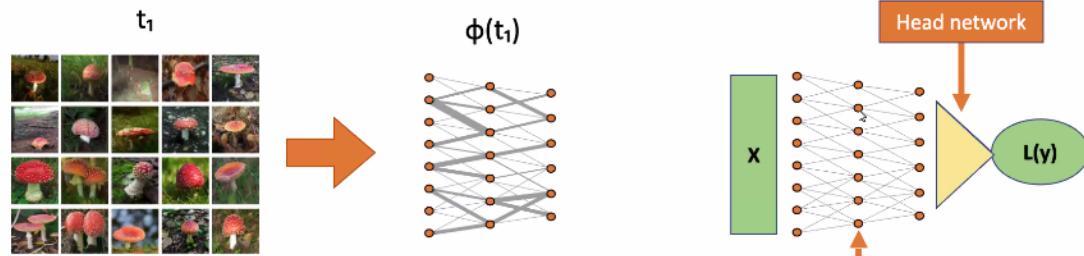
Optimal Λ satisfies:

$$\frac{\beta}{2N} \Lambda = F + \frac{\beta \lambda^2}{2N} I$$

↑
“Trivial Embedding”

14

Task embedding illustration

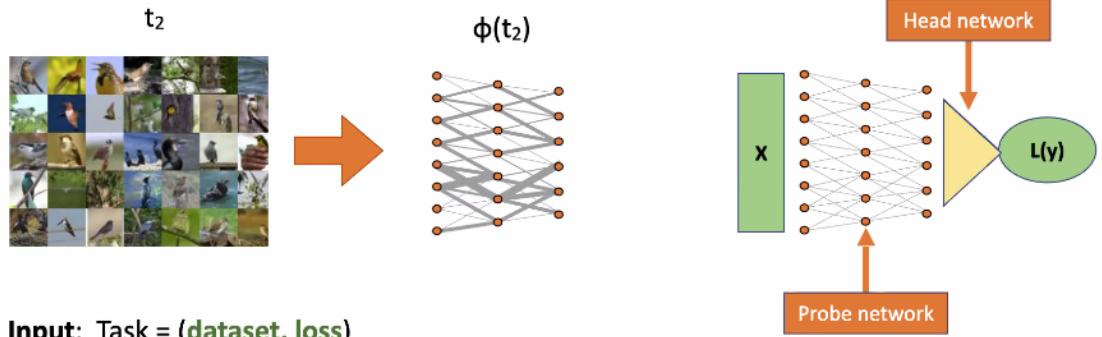


Input: Task = (**dataset, loss**)

- Initialize the **probe network** and the **head network** (e.g., **linear classifier**)
- Train the **head network** by minimizing the loss
- Compute the (approximate) FIM of the **probe network**

15

Task embedding illustration

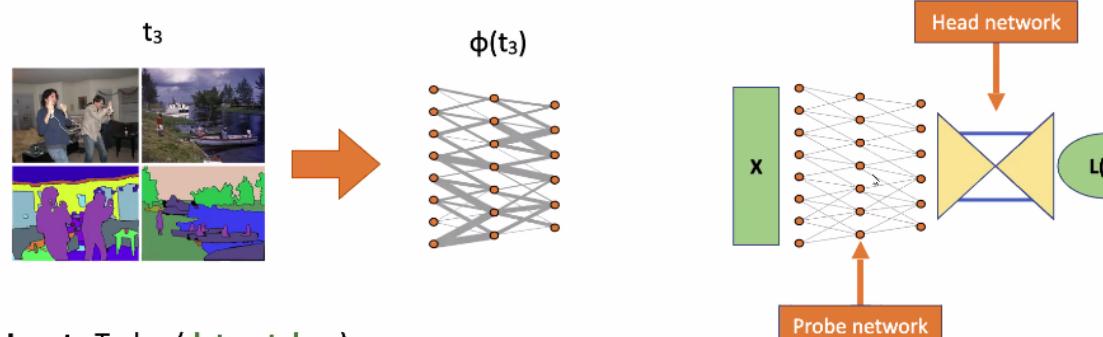


Input: Task = (dataset, loss)

1. Initialize the **probe network** and the **head network** (e.g., **linear classifier**)
2. Train the **head network** by minimizing the loss
3. Compute the (approximate) FIM of the **probe network**

16

Task embedding illustration

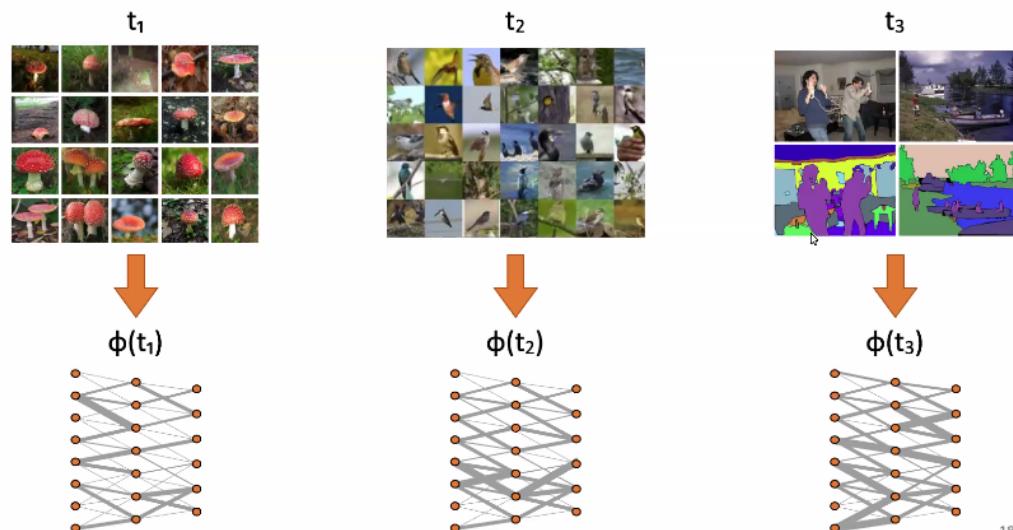


Input: Task = (dataset, loss)

1. Initialize the **probe network** and the **head network** (e.g., **UNet**)
2. Train the **head network** by minimizing the loss
3. Compute the (approximate) FIM of the **probe network**

17

Task embedding illustration



18

Properties of TASK2VEC embedding

Dataset

$$(x_i, y_i), i = 1 \dots n, \quad y_i \in \{0, 1\}$$

Classifier

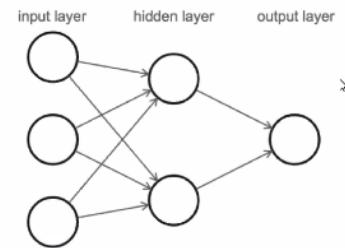
$$p_i = \sigma(w^T \phi(x_i))$$

FIM of the last layer (cross-entropy)

$$\frac{\partial \ell}{\partial w} = \frac{1}{N} \sum_i (y_i - p_i) \phi(x_i)$$

$$F_w = \frac{1}{N} \sum_i p_i (1 - p_i) \phi(x_i) \phi(x_i)^T$$

Two layer network



$$x \rightarrow \phi(x)$$

19

Properties of TASK2VEC embedding

Dataset

$$(x_i, y_i), i = 1 \dots n, \quad y_i \in \{0, 1\}$$

Classifier

$$p_i = \sigma(w^T \phi(x_i))$$

FIM of the last layer (cross-entropy)

$$\frac{\partial \ell}{\partial w} = \frac{1}{N} \sum_i (y_i - p_i) \phi(x_i)$$

$$F_w = \frac{1}{N} \sum_i p_i (1 - p_i) \phi(x_i) \phi(x_i)^T$$

1. **Invariance** to label space
2. Encodes task **difficulty**
3. Encodes task **domain**
4. Encodes **useful features** for the task

Representative “domain embedding”

$$D = \frac{1}{N} \sum_i \phi(x_i) \phi(x_i)^T$$

20

Distance measures on TASK2VEC embedding

Symmetric distance

$$d_{\text{sym}}(F_a, F_b) = d_{\cos}\left(\frac{F_a}{\|F_a\|}, \frac{F_b}{\|F_b\|}\right)$$

Task Segmentation

Asymmetric “distance”

$$d_{\text{asym}}(t_a \rightarrow t_b) = d_{\text{sym}}(t_a, t_b) - \alpha d_{\text{sym}}(t_a, t_0)$$

task embedding for the “trivial” task

21

Task Zoo

Tasks [1460]

- iNaturalist [207]
- CUB 200 [25]
- iMaterialist [228]
- DeepFashion [1000]



22

Task Zoo

Tasks [1460]

- iNaturalist [207]
- CUB 200 [25]
- iMaterialist [228]
- DeepFashion [1000]



24

Task Zoo

Tasks [1460]

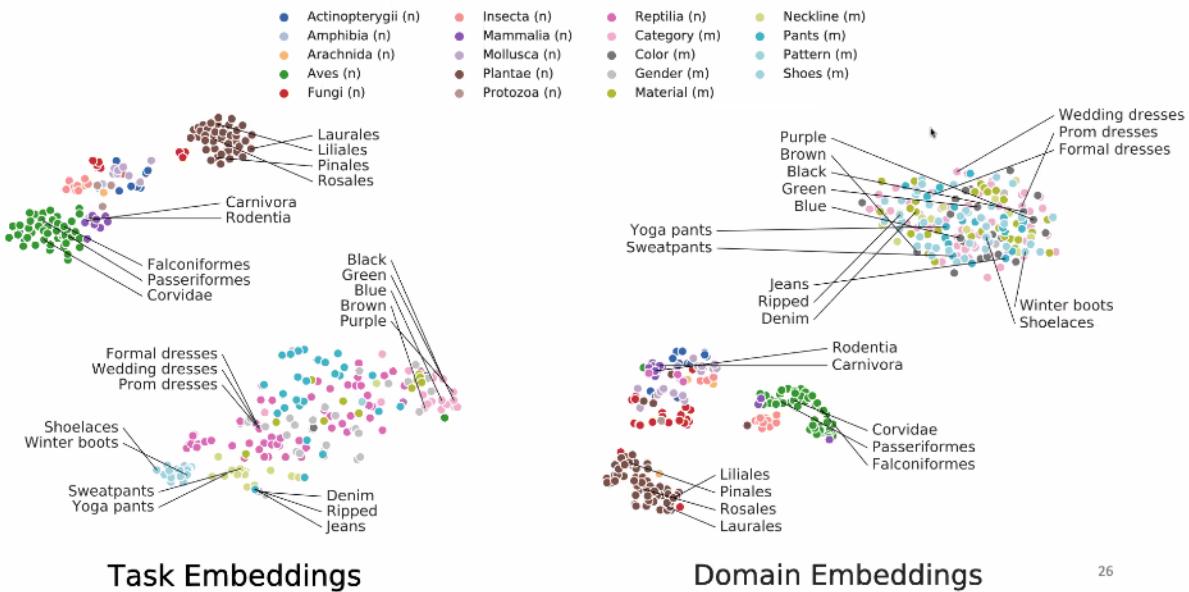
- iNaturalist [207]
- CUB 200 [25]
- iMaterialist [228]
- DeepFashion [1000]



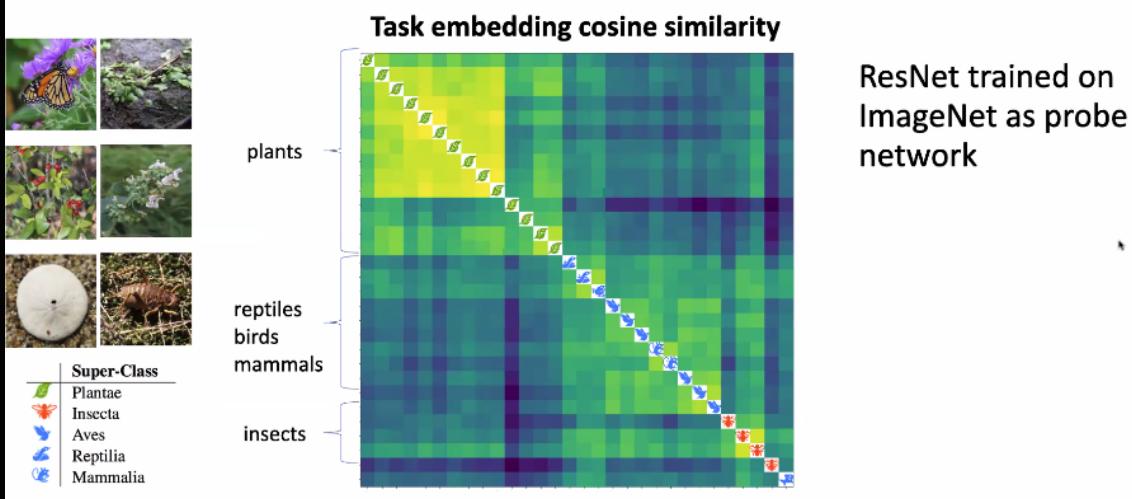
- Few tasks > 10K training samples but most have 100-1000 samples

25

Experiment: TASK2VEC vs DOMAIN2VEC

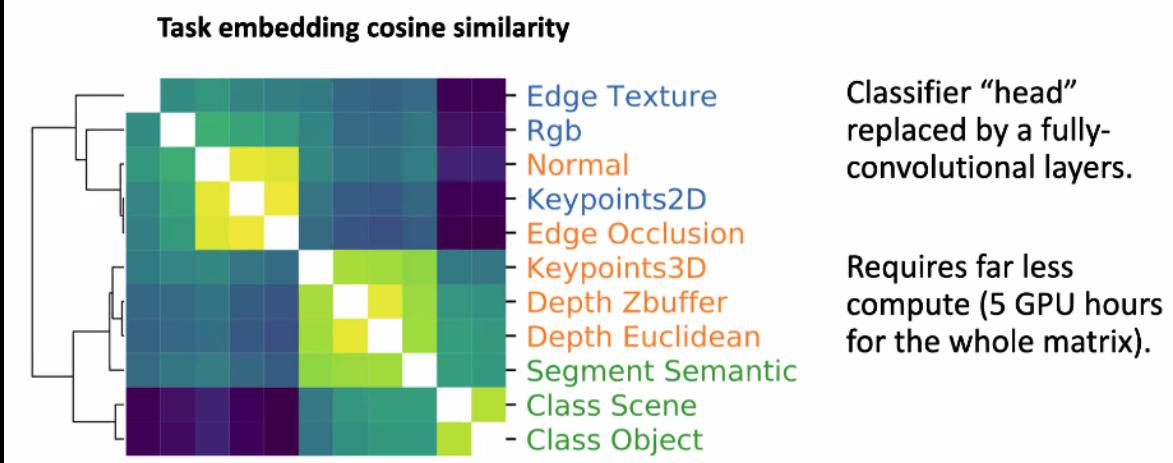


Experiment: TASK2VEC recapitulates iNaturalist taxonomy

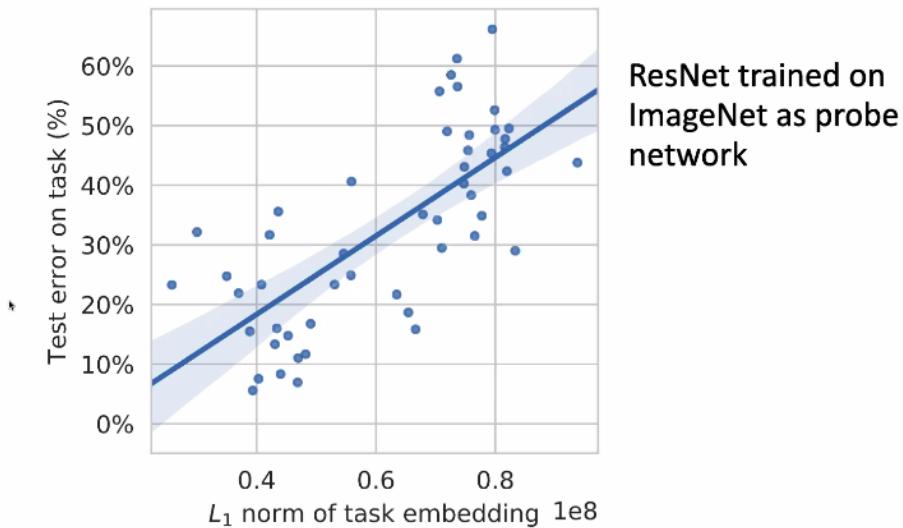


Experiment: TASK2VEC recovers “Taskonomy”

Taskonomy: Disentangling Task Transfer Learning, Amir Zamir, Alexander Sax, William Shen, Leonidas Guibas, Jitendra Malik, Silvio Savarese, CVPR 18

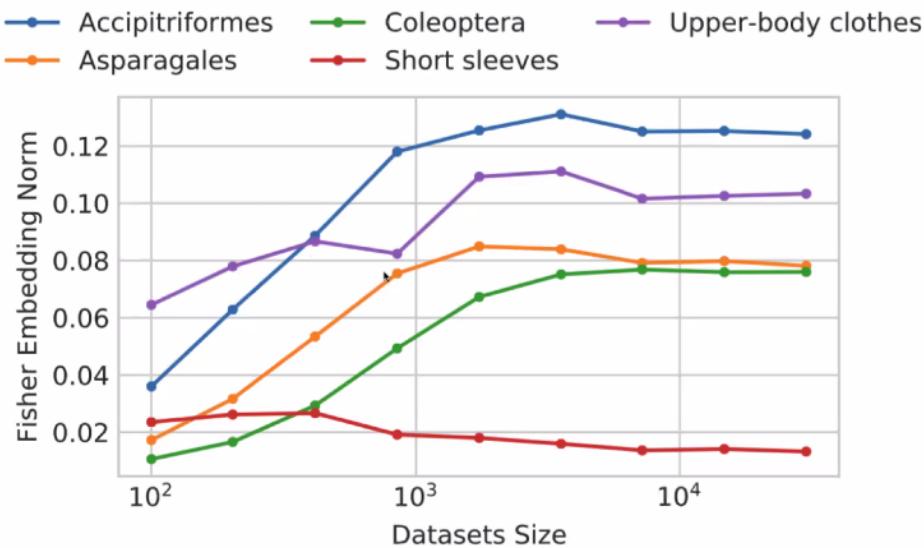


Experiment: Task2vec norm encodes task difficulty



29

Experiment: Task2vec norm encodes task difficulty



30

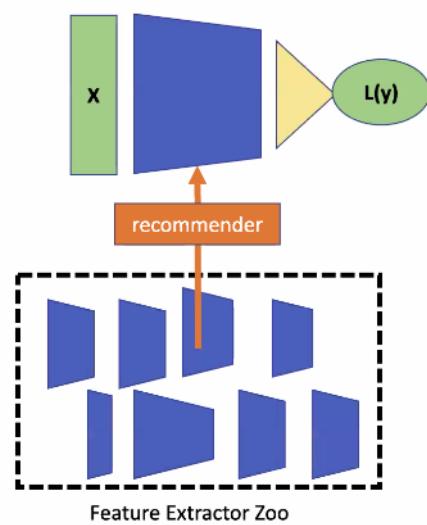
Task and Model Zoo

Tasks [1460]

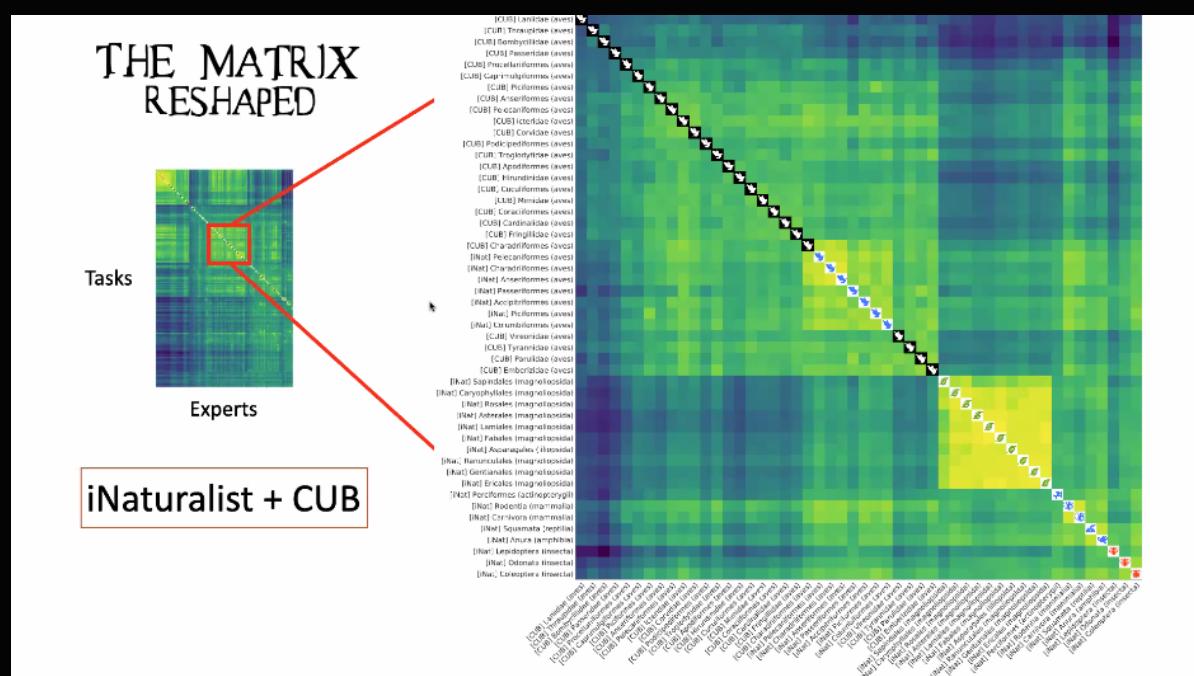
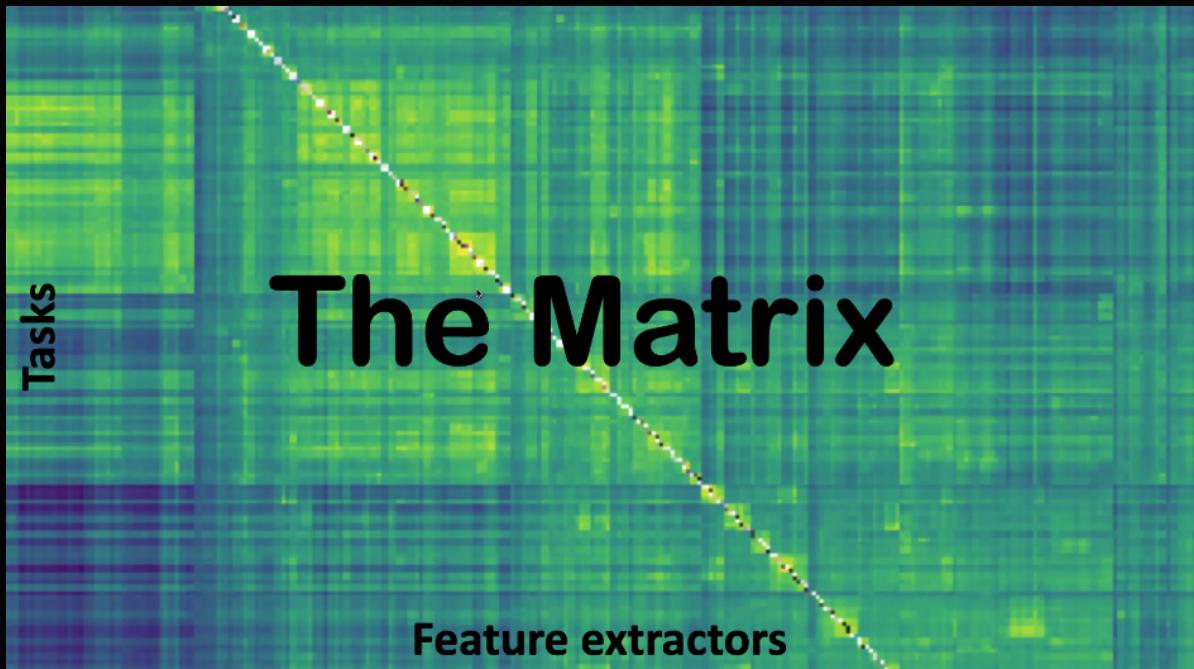
- iNaturalist [207]
- CUB 200 [25]
- iMaterialist [228]
- DeepFashion [1000]

Model Zoo [156 experts]

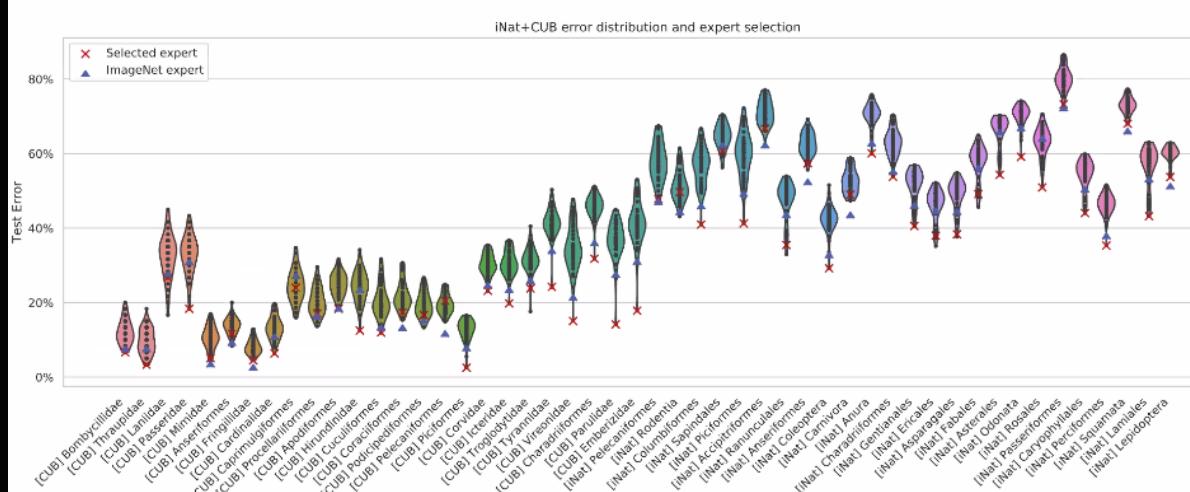
- ResNet-34 pertained on ImageNet
- Followed by fine-tuning on tasks with enough examples



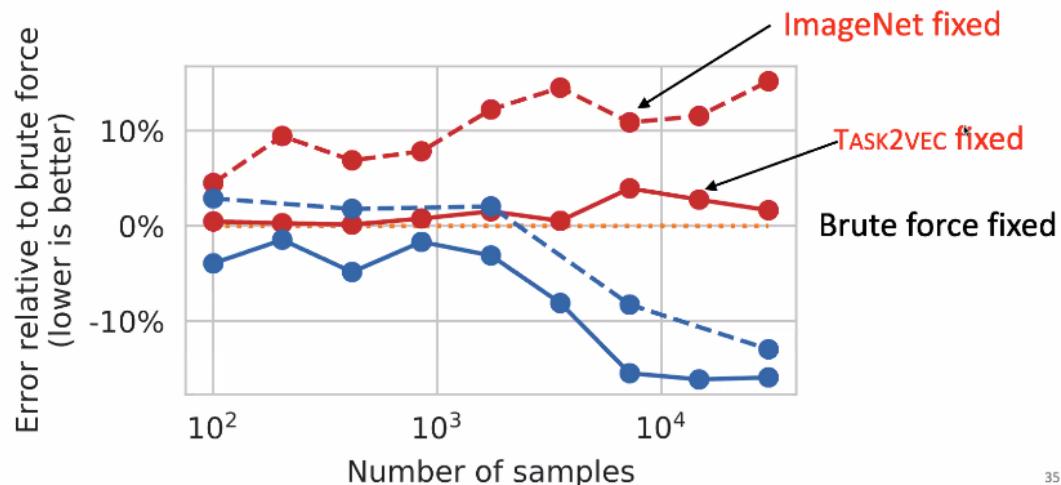
31



ImageNet expert is usually good, but on many tasks the best expert handily outperforms the ImageNet expert



Data efficiency of TASK2VEC



35

Talk outline

I. Task2Vec: vector representations of tasks [ICCV'19]

- Captures task similarity, task difficulty, etc.
- Useful for model recommendation



Alessandro Achille

II. When does self-supervision improve few-shot learning? [CVPR'20]

III. A realistic-evaluation of semi-supervised learning [CVPR'21]

- Self/semi—supervision is effective as long as the unlabeled images are from a similar domain
- Requires a notion of domain similarity



Jong-Chyi Su



Zezhou Cheng

Few-shot learning

Generalization to novel classes

- **Base classes:** many labeled examples
- **Novel classes:** few labeled examples
- **Base and novel classes** are disjoint



Train on
Base classes



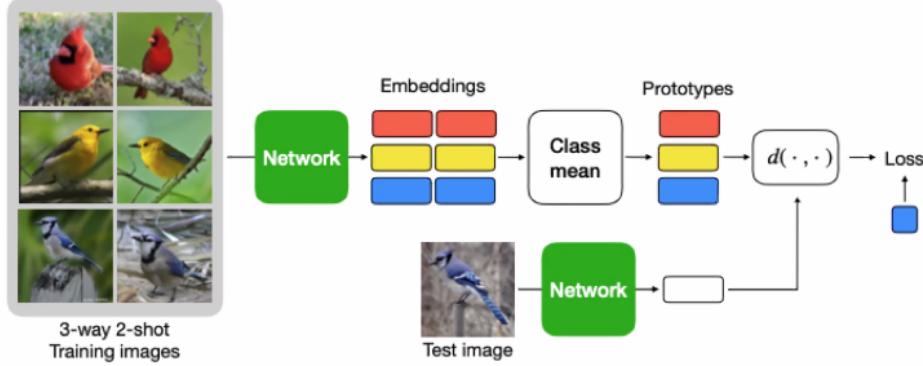
Test on
Novel classes

Many approaches based on meta-learning and metric learning

38

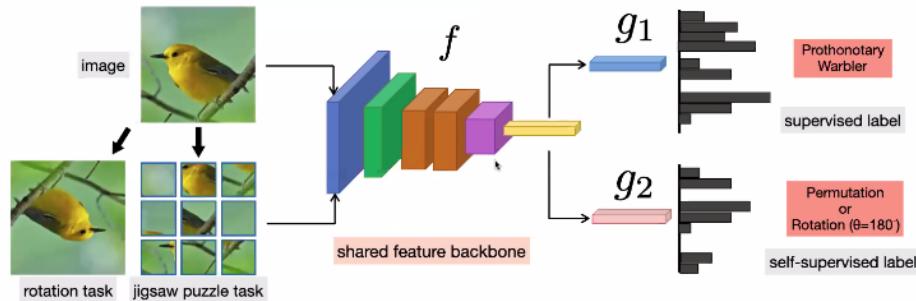
Prototypical networks [Snell et al., NeurIPS'17]

- Meta-learning with a nearest-mean classifier
- Sample tasks from the **base set**, minimize **error** on novel images



Self-supervision for few-shot learning

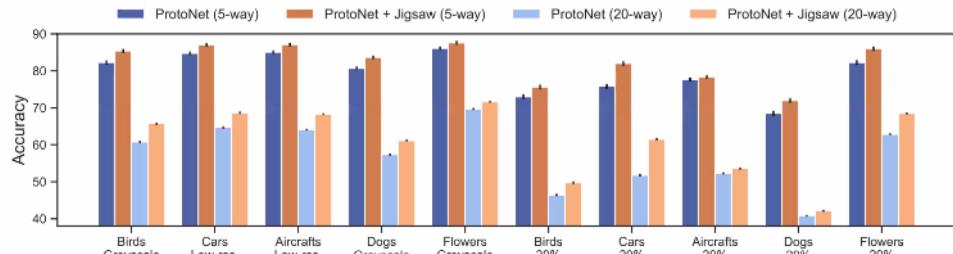
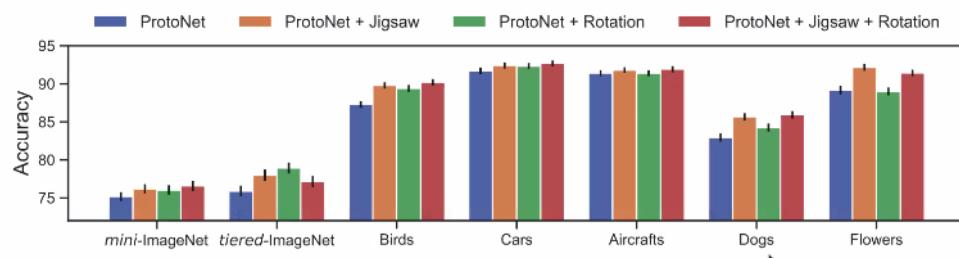
Jointly train the feature backbone with the few-shot and SSL objective



$$\min_{f, g_1, g_2} \mathbb{E}_{\substack{(x, y) \sim D_s, x' \sim D_{ss} \\ \text{supervised domain} \quad \text{self-supervised domain}}} [L_s(g_1 \circ f(x), y) + L_{ss}(g_2 \circ f(x'), y')]$$

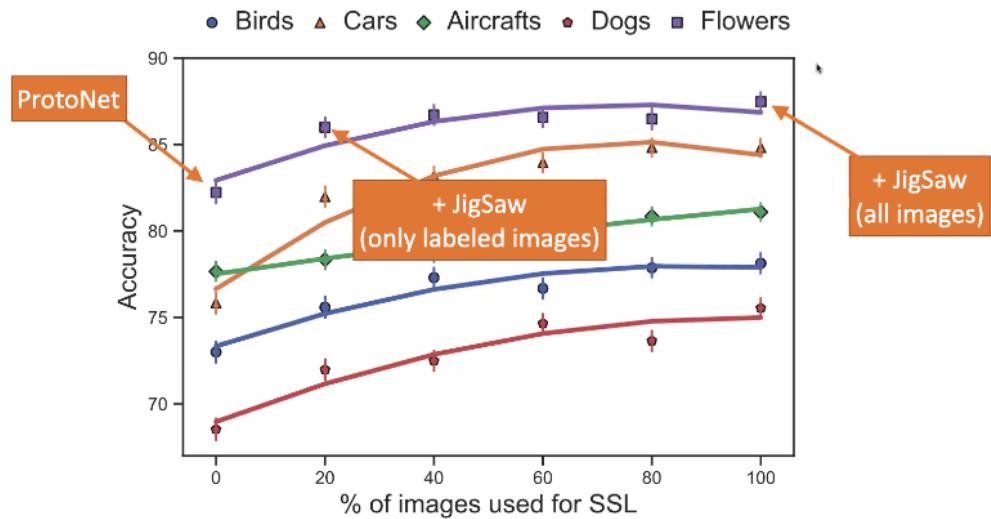
Also see: Gidaris et al., ICCV 2019

No extra data!
SSL helps even in the few-shot regime $D_s = D_{ss}$



extra data

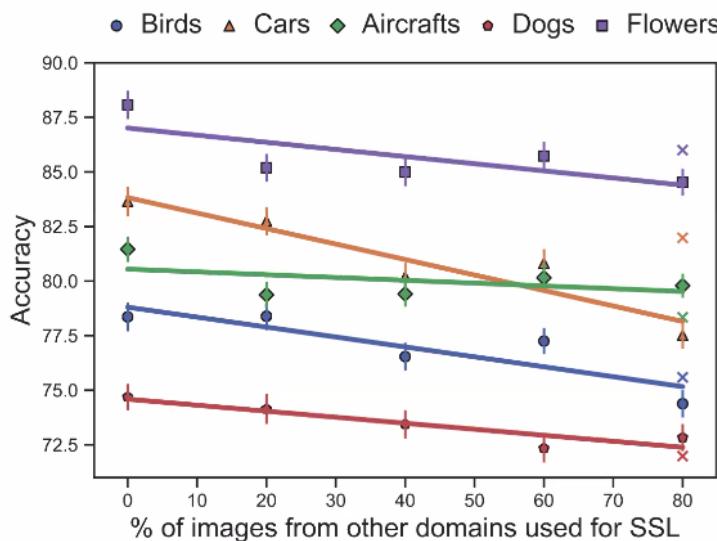
SSL helps even in the few-shot regime $D_s = D_{ss}$



42

extra data

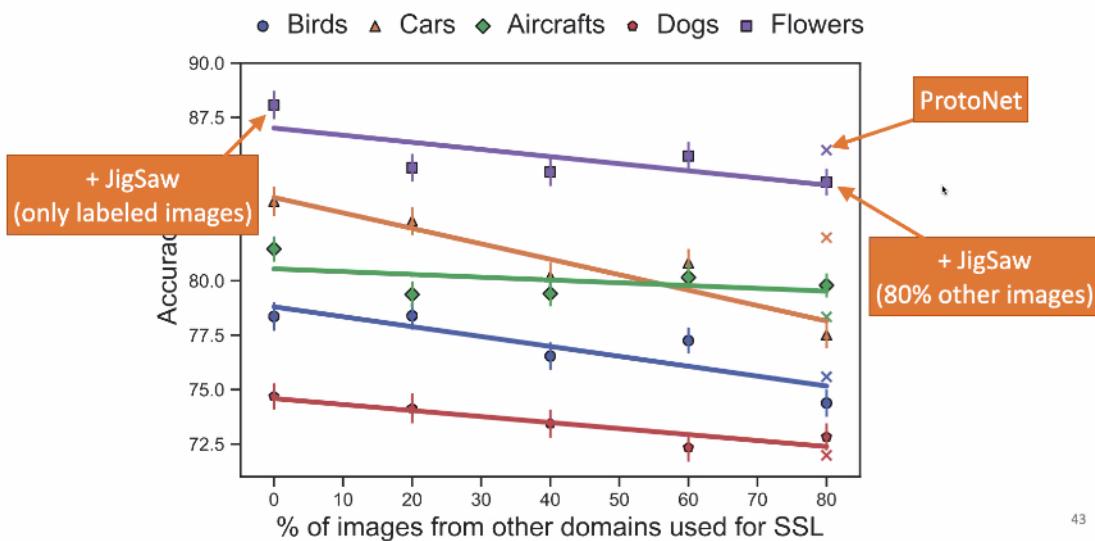
Domain shifts have a negative impact $D_s \neq D_{ss}$



43

extra data

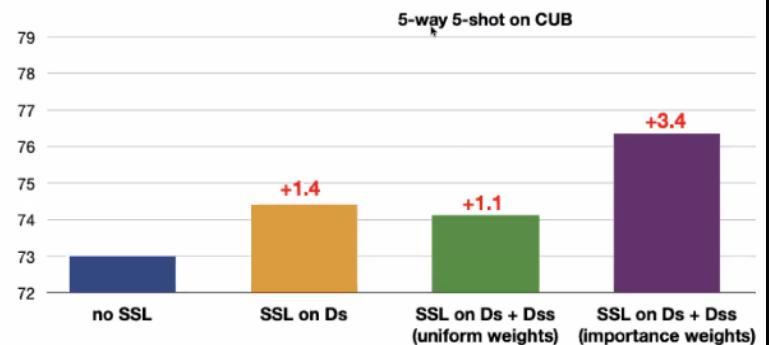
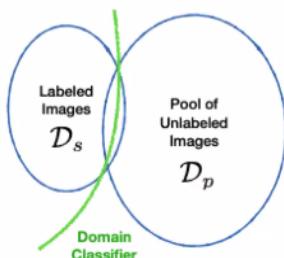
Domain shifts have a negative impact $D_s \neq D_{ss}$



43

Selecting images for self-supervision

- **Labeled images (\mathcal{D}_s):** 20% CUB dataset
- **Unlabeled images (\mathcal{D}_p):** OpenImages + iNaturalist (> 2M images)
- Select unlabeled images (\mathcal{D}_{ss}) for self-supervision



Talk outline

I. Task2Vec: Vector representations of tasks [ICCV'19]

- Captures task similarity, task difficulty, etc.
- Useful for model recommendation



Alessandro Achille

II. When does self-supervision improve few-shot learning? [CVPR'20]

III. A realistic-evaluation of semi-supervised learning [CVPR'21]

- Self/semi-supervision is effective as long as the unlabeled images are from a similar domain
- Requires a notion of domain similarity



Jong-Chyi Su

Zezhou Cheng

Benchmarking semi-supervised learning

A problem with existing benchmarks

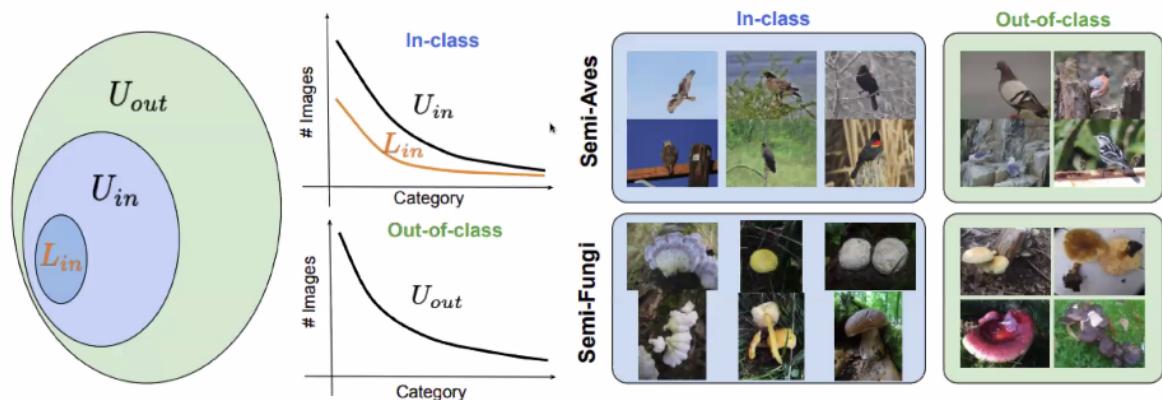
- Curated datasets: CIFAR, MNIST, SVHN, ImageNet
- Uniform class distribution
- Unlabeled data does not contain novel classes
- Role of transfer learning?

Motivated the FGVC7 and FGVC8 semi-supervised learning challenges

47

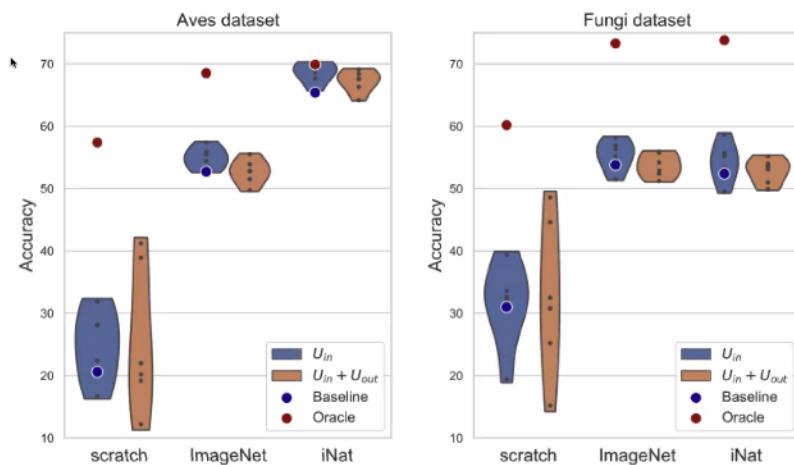
A realistic benchmark

Semi-Aves and Semi-Fungi datasets (CVPR'20)



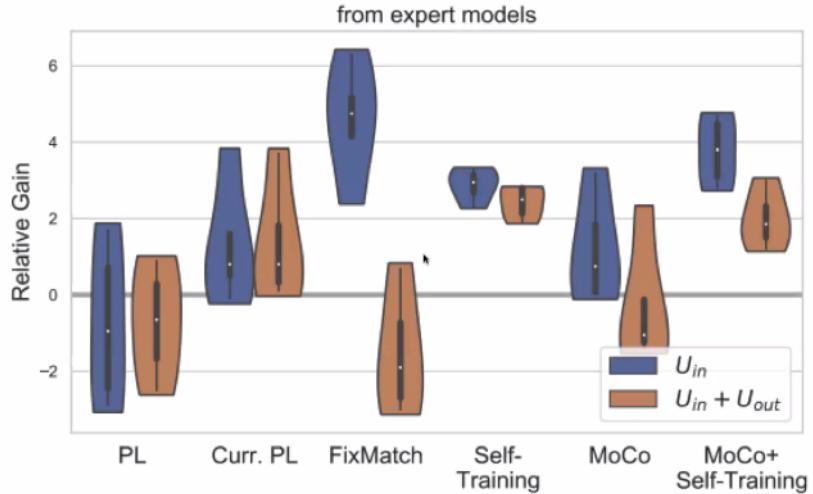
48

Takeaway #1: Transfer learning is effective, levels the playing field for semi-supervised learning



49

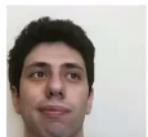
Takeaway #2: Methods are not robust to domain shifts (e.g., presence of novel classes)



50

Thank you!

- I. **Task2Vec**: Vector representations of tasks [ICCV'19]
- II. **When does self-supervision improve few-shot learning?** [CVPR'20]
- III. **A realistic-evaluation of semi-supervised learning** [CVPR'21]



+ Many others



UMassAmherst

College of Information and Computer Sciences



51