

# Visual Learning Beyond Natural Images

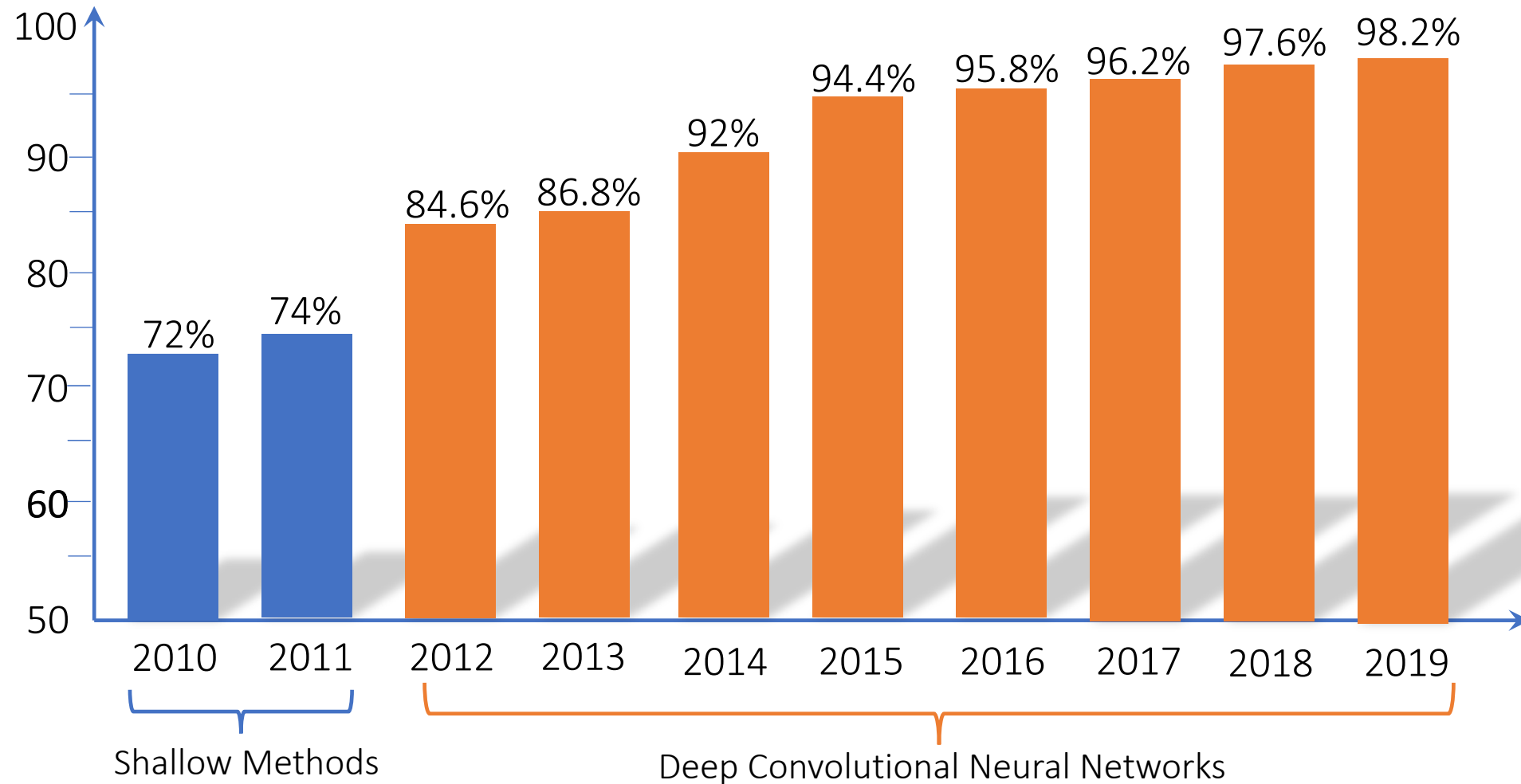
Rogério Schmidt Feris

MIT-IBM Watson AI Lab

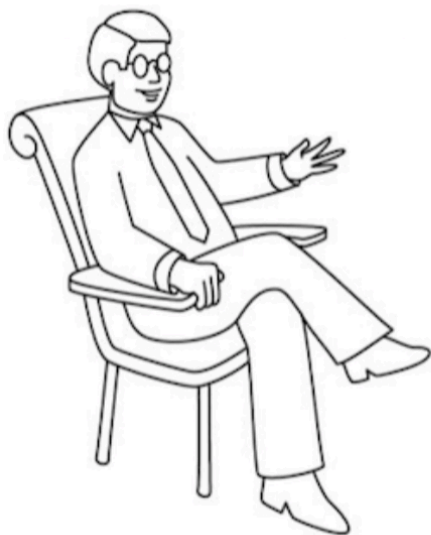
<http://rogerioferis.com>



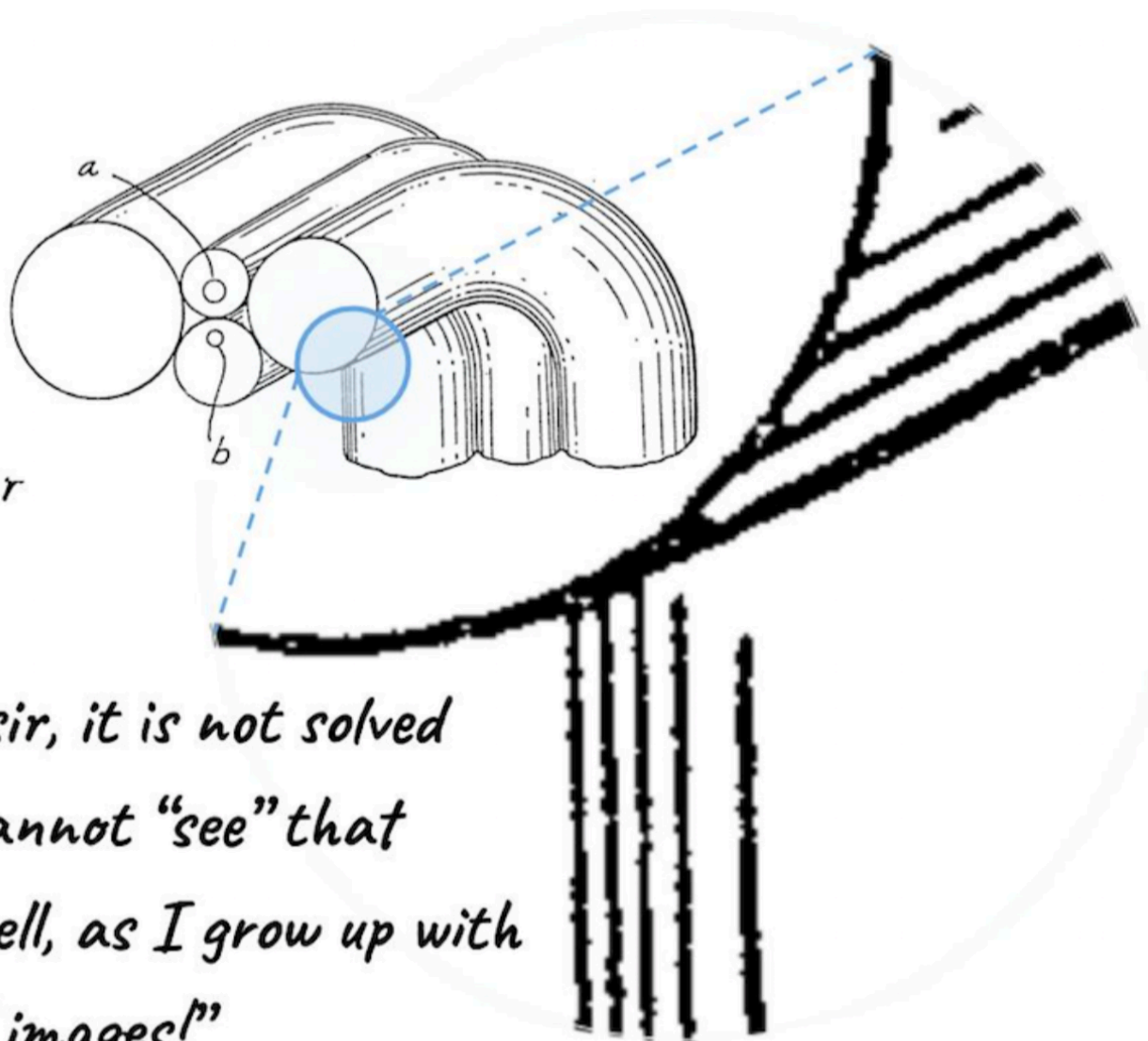
# ImageNet Classification (top-5 accuracy)



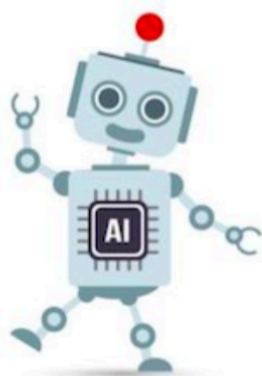
# However...



*"Bob, can you tell me whether this design has been patented or not?"*



*"sorry, sir, it is not solved yet; I cannot "see" that image well, as I grow up with natural images!"*

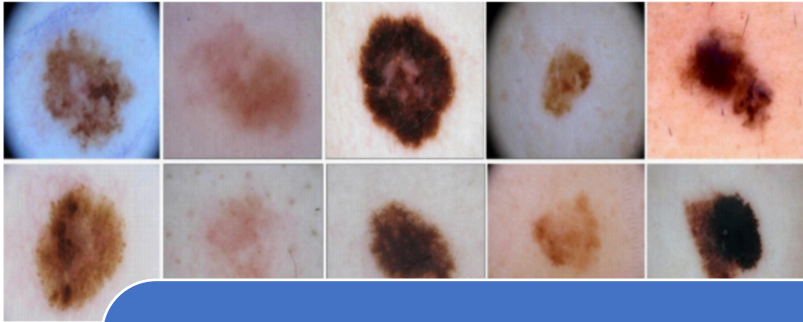






# Visual Learning Beyond Natural Images

Healthcare



Sattelite



Document Analysis

R.J.Reynolds Tobacco Company

AUTHORIZATION REQUEST *mk-137*  
*SC 4 - JCB*

PROJECT: CAMEL Theme Promotion Contract Amendment DATE PREPARED: 7/8/85 AR NO.: 85-479  
ORIGINATOR: C. L. Sharp DEPT.: 853

APPROVAL REQUEST SUMMARY

This requests Management approval to amend the CAMEL Theme Promotion contract with Glendinning Associates as follows:

- Fully develop and submit to RJRT Interim Tactical Plans for the following promotion concepts previously developed by Glendinning Associates:
  - Live the Adventure Sweepstakes
  - Keep Your Machine in Winning Shape Free Premium Mail-In

Problem: Limited Labeled Data



... and Much More !

COMPLETE IF LEASE OR OTHER CONTINUING COMMITMENT IS INVOLVED

Commitment \$ \_\_\_\_\_ Per \_\_\_\_\_ For \_\_\_\_\_ Years  
Total Commitment \$ \_\_\_\_\_ Minimum Commitment \$ \_\_\_\_\_ Time \_\_\_\_\_  
Return on Invest. \_\_\_\_\_ %

EFO IMPACT 19 85 19 \_\_\_\_\_ 19 \_\_\_\_\_ Annual Average  
Profit/(Loss) (\$25,989) \_\_\_\_\_  
Budget/Plan Change  Yes  No

REVIEWED BY:				APPROVALS (Originator Enters Initials of Required Approvals)			
Dept.	Initials	Date	Action Completed	Initials	Signature	Date	
Prom.	BV	7/8	✓	MLO	<i>W. Sharp</i>	7/8	
Prom.	CLS	7/8	✓				
Prom.	TBO	7/9	✓				
Gen. Sv.	CTB	7/10	✓				EFO
Corp.	DRS	7/11	✓				Operations
Law	MT	7/11	✓				Fin. & Admin.
Risk Mgt	M.B.	7/11	✓				Fin. & CEO

Person Responsible for Implementing: C. L. Sharp Project ID No. \_\_\_\_\_

51316 1442



# A Broader Study of Cross-Domain Few-Shot Learning

Yunhui Guo, Noel C Codella, Leonid Karlinsky, James V Codella,  
John R Smith, Kate Saenko, Tajana Rosing, Rogerio Feris

ECCV 2020



# Few-Shot Learning Problem

## Meta-Learning

Meta-learning: Learn how to learn with few examples in training tasks. Performed via many trials of k-shot n-way tasks, while optimizing.

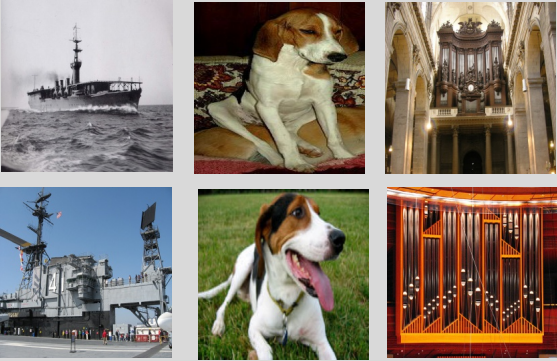
## Meta-Testing

Evaluate method on novel tasks and measure accuracy over many trials.

Support set:

“Episode”

“k-shot”

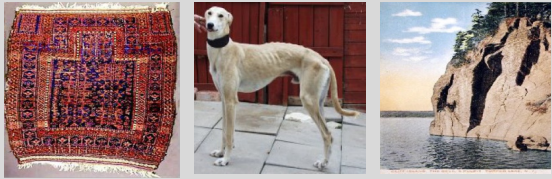
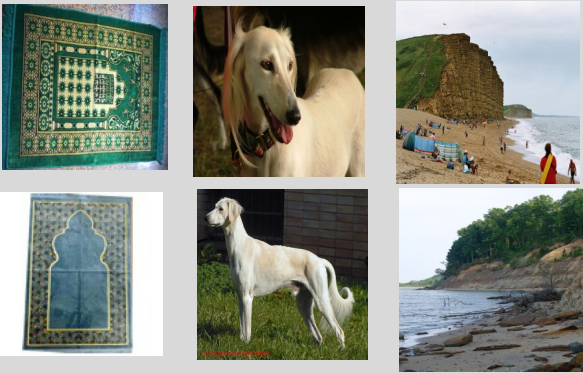


“n-way”

Query set:



...



# What happens when data is different from expected domain?

k-way

## Meta-Learning (TRAIN)

"House Finch"



"Gazelle Hound"



"Triceratops"



n-shot

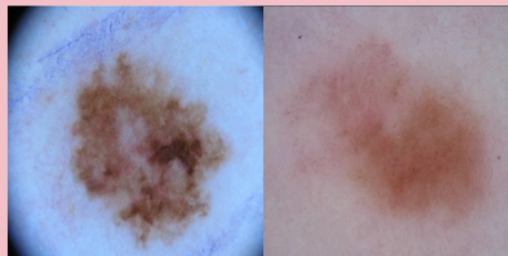
## "Novel" Category (TEST – SAME DOMAIN)

"Golden Retriever"



## "Novel" Category (TEST – DIFF. DOMAIN)

"Melanoma"



Current research explores few-shot categories with exceptionally high similarity to training data!

In practice, few-shot categories can wildly differ to training data, causing all existing techniques to break-down and perform worse than baseline simple methods.





# tieredImageNet was first attempt to address

The **miniImageNet** dataset [33] is a standard benchmark for few-shot image classification benchmark, consisting of 100 randomly chosen classes from ILSVRC-2012 [24]. These classes are randomly split into 64, 16 and 20 classes for meta-training, meta-validation, and meta-testing respectively. Each class contains 600 images of size  $84 \times 84$ . Since the class splits were not released in the original publication [33], we use the commonly-used split proposed in [22].

The **tieredImageNet** benchmark [23] is a larger subset of ILSVRC-2012 [24], composed of 608 classes grouped into 34 high-level categories. These are divided into 20 categories for meta-training, 6 categories for meta-validation, and 8 categories for meta-testing. This corresponds to 351, 97 and 160 classes for meta-training, meta-validation, and meta-testing respectively. This dataset aims to minimize the semantic similarity between the splits. All images are of size  $84 \times 84$ .

But these are all still categories within the domain of natural images!





# Significant progress on minImageNet / tieredImageNet

(numbers below already outdated)

model	backbone	minImageNet 5-way		tieredImageNet 5-way	
		1-shot	5-shot	1-shot	5-shot
Meta-Learning LSTM* [22]	64-64-64-64	43.44 ± 0.77	60.60 ± 0.71	-	-
Matching Networks* [33]	64-64-64-64	43.56 ± 0.84	55.31 ± 0.73	-	-
MAML [8]	32-32-32-32	48.70 ± 1.84	63.11 ± 0.92	51.67 ± 1.81	70.30 ± 1.75
Prototypical Networks* <sup>†</sup> [28]	64-64-64-64	49.42 ± 0.78	68.20 ± 0.66	53.31 ± 0.89	72.69 ± 0.74
Relation Networks* [29]	64-96-128-256	50.44 ± 0.82	65.32 ± 0.70	54.48 ± 0.93	71.32 ± 0.78
R2D2 [3]	96-192-384-512	51.2 ± 0.6	68.8 ± 0.1	-	-
Transductive Prop Nets [14]	64-64-64-64	55.51 ± 0.86	69.86 ± 0.65	59.91 ± 0.94	73.30 ± 0.75
SNAIL [18]	ResNet-12	55.71 ± 0.99	68.88 ± 0.92	-	-
Dynamic Few-shot [10]	64-64-128-128	56.20 ± 0.86	73.00 ± 0.64	-	-
AdaResNet [19]	ResNet-12	56.88 ± 0.62	71.94 ± 0.57	-	-
TADAM [20]	ResNet-12	58.50 ± 0.30	76.70 ± 0.30	-	-
Activation to Parameter <sup>†</sup> [21]	WRN-28-10	59.60 ± 0.41	73.74 ± 0.19	-	-
LEO <sup>†</sup> [25]	WRN-28-10	61.76 ± 0.08	77.59 ± 0.12	<b>66.33 ± 0.05</b>	<b>81.44 ± 0.09</b>
MetaOptNet-RR (ours)	ResNet-12	61.41 ± 0.61	77.88 ± 0.46	<b>65.36 ± 0.71</b>	<b>81.34 ± 0.52</b>
MetaOptNet-SVM (ours)	ResNet-12	62.64 ± 0.61	78.63 ± 0.46	<b>65.99 ± 0.72</b>	<b>81.56 ± 0.53</b>
MetaOptNet-SVM-trainval (ours) <sup>†</sup>	ResNet-12	<b>64.09 ± 0.62</b>	<b>80.00 ± 0.45</b>	<b>65.81 ± 0.74</b>	<b>81.75 ± 0.53</b>

# Let's get more realistic by gradually leaving domain: Proposed Cross-Domain Evaluation Benchmark

Source Domain:



**ImageNet:**  
Perspective  
Natural Images  
Color

**Target Domains:**

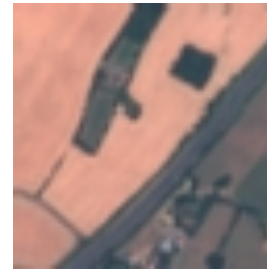
(Disjoint Label Spaces)



Decreasing Similarity to ImageNet



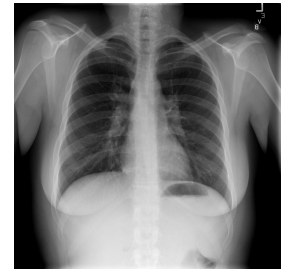
**CropDisease:**  
Perspective  
Natural Images  
Color



**EuroSAT:**  
No Perspective  
Natural Images  
Color



**ISIC:**  
No Perspective  
Medical Images  
Color



**ChestX:**  
No Perspective  
Medical Images  
Grayscale

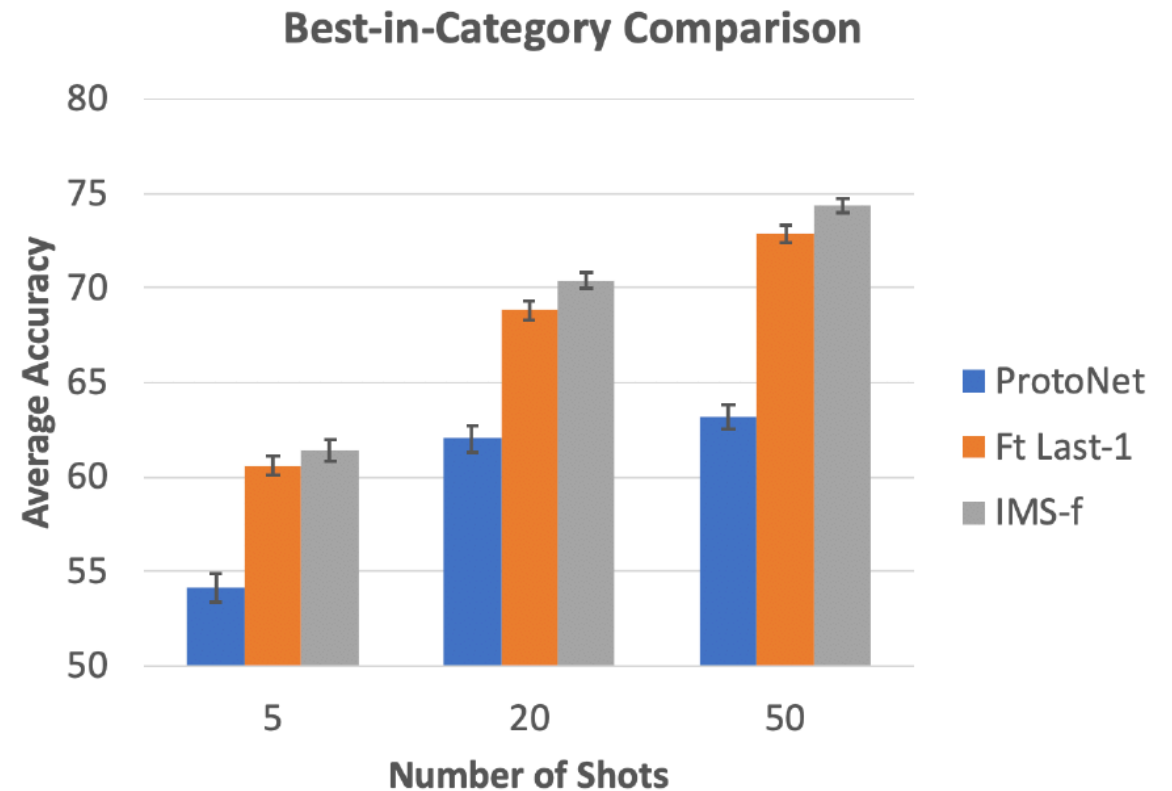
<https://www.learning-with-limited-labels.com/challenge>



# Meta-learning doesn't generalize well. Fine-tuning is better. Ensembles are even better.

- We evaluate performance of current meta-learning, cross-domain few-shot learning, fine-tuning, and ensemble methods.
- Meta-learning approaches perform worst, even methods tailored for cross-domain.
- Fine-tuning is better.
- Ensemble approaches are best.

<https://arxiv.org/pdf/1912.07200.pdf>



**Fig. 4:** Best meta-learning, single model, and multi-model transfer learning.



# Details of results on each dataset

Best Meta-learning approach

Prior state-of-art

Methods	ChestX			ISIC		
	5-way 5-shot	5-way 20-shot	5-way 50-shot	5-way 5-shot	5-way 20-shot	5-way 50-shot
<i>MatchingNet</i>	22.40% $\pm$ 0.7%	23.61% $\pm$ 0.86%	22.12% $\pm$ 0.88%	36.74% $\pm$ 0.53%	45.72% $\pm$ 0.53%	54.58% $\pm$ 0.65%
<i>MatchingNet+FWT</i>	21.26% $\pm$ 0.31%	23.23% $\pm$ 0.37%	23.01% $\pm$ 0.34%	30.40% $\pm$ 0.48%	32.01% $\pm$ 0.48%	33.17% $\pm$ 0.43%
<i>MAML</i>	23.48% $\pm$ 0.96%	27.53% $\pm$ 0.43%	-	40.13% $\pm$ 0.58%	52.36% $\pm$ 0.57%	-
<i>ProtoNet</i>	24.05% $\pm$ 1.01%	28.21% $\pm$ 1.15%	29.32% $\pm$ 1.12%	39.57% $\pm$ 0.57%	49.50% $\pm$ 0.55%	51.99% $\pm$ 0.52%
<i>ProtoNet+FWT</i>	23.77% $\pm$ 0.42%	26.87% $\pm$ 0.43%	30.12% $\pm$ 0.46%	38.87% $\pm$ 0.52%	43.78% $\pm$ 0.47%	49.84% $\pm$ 0.51%
<i>RelationNet</i>	22.96% $\pm$ 0.88%	26.63% $\pm$ 0.92%	28.45% $\pm$ 1.20%	39.41% $\pm$ 0.58%	41.77% $\pm$ 0.49%	49.32% $\pm$ 0.51%
<i>RelationNet+FWT</i>	22.74% $\pm$ 0.40%	26.75% $\pm$ 0.41%	27.56% $\pm$ 0.40%	35.54% $\pm$ 0.55%	43.31% $\pm$ 0.51%	46.38% $\pm$ 0.53%
<i>MetaOpt</i>	22.53% $\pm$ 0.91%	25.53% $\pm$ 1.02%	29.35% $\pm$ 0.99%	36.28% $\pm$ 0.50%	49.42% $\pm$ 0.60%	54.80% $\pm$ 0.54%

Methods	EuroSAT			CropDiseases		
	5-way 5-shot	5-way 20-shot	5-way 50-shot	5-way 5-shot	5-way 20-shot	5-way 50-shot
<i>MatchingNet</i>	64.45% $\pm$ 0.63%	77.10% $\pm$ 0.57%	54.44% $\pm$ 0.67%	66.39% $\pm$ 0.78%	76.38% $\pm$ 0.67%	58.53% $\pm$ 0.73%
<i>MatchingNet+FWT</i>	56.04% $\pm$ 0.65%	63.38% $\pm$ 0.69%	62.75% $\pm$ 0.76%	62.74% $\pm$ 0.90%	74.90% $\pm$ 0.71%	75.68% $\pm$ 0.78%
<i>MAML</i>	71.70% $\pm$ 0.72%	81.95% $\pm$ 0.55%	-	78.05% $\pm$ 0.68%	89.75% $\pm$ 0.42%	-
<i>ProtoNet</i>	73.29% $\pm$ 0.71%	82.27% $\pm$ 0.57%	80.48% $\pm$ 0.57%	79.72% $\pm$ 0.67%	88.15% $\pm$ 0.51%	90.81% $\pm$ 0.43%
<i>ProtoNet+FWT</i>	67.34% $\pm$ 0.76%	75.74% $\pm$ 0.70%	78.64% $\pm$ 0.57%	72.72% $\pm$ 0.70%	85.82% $\pm$ 0.51%	87.17% $\pm$ 0.50%
<i>RelationNet</i>	61.31% $\pm$ 0.72%	74.43% $\pm$ 0.66%	74.91% $\pm$ 0.58%	68.99% $\pm$ 0.75%	80.45% $\pm$ 0.64%	85.08% $\pm$ 0.53%
<i>RelationNet+FWT</i>	61.16% $\pm$ 0.70%	69.40% $\pm$ 0.64%	73.84% $\pm$ 0.60%	64.91% $\pm$ 0.79%	78.43% $\pm$ 0.59%	81.14% $\pm$ 0.56%
<i>MetaOpt</i>	64.44% $\pm$ 0.73%	79.19% $\pm$ 0.62%	83.62% $\pm$ 0.58%	68.41% $\pm$ 0.73%	82.89% $\pm$ 0.54%	91.76% $\pm$ 0.38%

**Table 1:** The results of meta-learning methods on the proposed benchmark.





# Details of results on each dataset

## Best transfer-learning approach

Methods	ChestX			ISIC		
	5-way 5-shot	5-way 20-shot	5-way 50-shot	5-way 5-shot	5-way 20-shot	5-way 50-shot
<i>Random</i>	21.80% $\pm$ 1.03%	25.69% $\pm$ 0.95%	26.19% $\pm$ 0.94%	37.91% $\pm$ 1.39%	47.24% $\pm$ 1.50%	50.85% $\pm$ 1.37%
<i>Fixed</i>	25.35% $\pm$ 0.96%	30.83% $\pm$ 1.05%	36.04% $\pm$ 0.46%	43.56% $\pm$ 0.60%	52.78% $\pm$ 0.58%	57.34% $\pm$ 0.56%
<i>Ft All</i>	25.97% $\pm$ 0.41%	31.32% $\pm$ 0.45%	35.49% $\pm$ 0.45%	48.11% $\pm$ 0.64%	59.31% $\pm$ 0.48%	66.48% $\pm$ 0.56%
<i>Ft Last-1</i>	25.96% $\pm$ 0.46%	31.63% $\pm$ 0.49%	37.03% $\pm$ 0.50%	47.20% $\pm$ 0.45%	59.95% $\pm$ 0.45%	65.04% $\pm$ 0.47%
<i>Ft Last-2</i>	26.79% $\pm$ 0.59%	30.95% $\pm$ 0.61%	36.24% $\pm$ 0.62%	47.64% $\pm$ 0.44%	59.87% $\pm$ 0.35%	66.07% $\pm$ 0.45%
<i>Ft Last-3</i>	25.17% $\pm$ 0.56%	30.92% $\pm$ 0.89%	37.27% $\pm$ 0.64%	48.05% $\pm$ 0.55%	60.20% $\pm$ 0.33%	66.21% $\pm$ 0.52%
<i>Transductive Ft</i>	26.09% $\pm$ 0.96%	31.01% $\pm$ 0.59%	36.79% $\pm$ 0.53%	49.68% $\pm$ 0.36%	61.09% $\pm$ 0.44%	67.20% $\pm$ 0.59%

Methods	EuroSAT			CropDiseases		
	5-way 5-shot	5-way 20-shot	5-way 50-shot	5-way 5-shot	5-way 20-shot	5-way 50-shot
<i>Random</i>	58.00% $\pm$ 2.01%	68.93% $\pm$ 1.47%	71.65% $\pm$ 1.47%	69.68% $\pm$ 1.72%	83.41% $\pm$ 1.25%	86.56% $\pm$ 1.42%
<i>Fixed</i>	75.69% $\pm$ 0.66%	84.13% $\pm$ 0.52%	86.62% $\pm$ 0.47%	87.48% $\pm$ 0.58%	94.45% $\pm$ 0.36%	96.62% $\pm$ 0.25%
<i>Ft All</i>	79.08% $\pm$ 0.61%	87.64% $\pm$ 0.47%	90.89% $\pm$ 0.36%	89.25% $\pm$ 0.51%	95.51% $\pm$ 0.31%	97.68% $\pm$ 0.21%
<i>Ft Last-1</i>	80.45% $\pm$ 0.54%	87.92% $\pm$ 0.44%	91.41% $\pm$ 0.46%	88.72% $\pm$ 0.53%	95.76% $\pm$ 0.65%	97.87% $\pm$ 0.48%
<i>Ft Last-2</i>	79.57% $\pm$ 0.51%	87.67% $\pm$ 0.46%	90.93% $\pm$ 0.45%	88.07% $\pm$ 0.56%	95.68% $\pm$ 0.76%	97.64% $\pm$ 0.59%
<i>Ft Last-3</i>	78.04% $\pm$ 0.77%	87.52% $\pm$ 0.53%	90.83% $\pm$ 0.42%	89.11% $\pm$ 0.47%	95.31% $\pm$ 0.7%	97.45% $\pm$ 0.46%
<i>Transductive Ft</i>	81.76% $\pm$ 0.48%	87.97% $\pm$ 0.42%	92.00% $\pm$ 0.56%	90.64% $\pm$ 0.54%	95.91% $\pm$ 0.72%	97.48% $\pm$ 0.56%

**Table 2:** The results of different variants of single model fine-tuning on the proposed benchmark.





# Intelligent selective ensembles outperform naïve ensembles.

Best ensemble approach

Methods	ChestX			ISIC		
	5-way 5-shot	5-way 20-shot	5-way 50-shot	5-way 5-shot	5-way 20-shot	5-way 50-shot
<i>All embeddings</i>	26.74% $\pm$ 0.42%	32.77% $\pm$ 0.47%	38.07% $\pm$ 0.50%	46.86% $\pm$ 0.60%	58.57% $\pm$ 0.59%	66.04% $\pm$ 0.56%
<i>IMS-f</i>	25.50% $\pm$ 0.45%	31.49% $\pm$ 0.47%	36.40% $\pm$ 0.50%	45.84% $\pm$ 0.62%	61.50% $\pm$ 0.58%	68.64% $\pm$ 0.53%

Methods	EuroSAT			CropDiseases		
	5-way 5-shot	5-way 20-shot	5-way 50-shot	5-way 5-shot	5-way 20-shot	5-way 50-shot
<i>All embeddings</i>	81.29% $\pm$ 0.62%	89.90% $\pm$ 0.41%	92.76% $\pm$ 0.34%	90.82% $\pm$ 0.48%	96.64% $\pm$ 0.25%	98.14% $\pm$ 0.18%
<i>IMS-f</i>	83.56% $\pm$ 0.59%	91.22% $\pm$ 0.38%	93.85% $\pm$ 0.30%	90.66% $\pm$ 0.48%	97.18% $\pm$ 0.24%	98.43% $\pm$ 0.16%

**Table 4:** The results of using all embeddings, and the *Incremental Multi-model Selection* (IMS-f) based on fine-tuned pre-trained models on the proposed benchmark.



# Similar Conclusions Hold for Document Analysis

Experiments with the RVL-CDIP dataset (document classification)

	<b>5-way 5-shot</b>	<b>5-way 20-shot</b>	<b>5-way 1-shot</b>	<b>5-way 50-shot</b>
MatchinetNet	42.18% +- 0.72%	43.41% +- 0.64%	29.45% +- 0.59%	27.10% +- 0.48%
ProtoNet	49.92% +- 0.81%	55.74% +- 0.74%	35.27% +- 0.72%	57.20% +- 0.75%
MetaOpt	41.11% +- 0.72%	54.60% +- 0.75%	33.86% +- 0.71%	62.97% +- 0.72%
RelationNet	44.09% +- 0.71%	52.39% +- 0.70%	35.66% +- 0.79%	55.57% +- 0.66%
MAML	34.48% +- 0.69%	36.49% +- 0.67%	36.23% +- 0.80%	-
Fixed	51.31% +- 0.78%	61.82% +- 0.72%	35.03% +- 0.74%	66.30% +- 0.69%
Fine-tune	55.93% +- 0.79%	64.78% +- 0.75%	37.01% +- 0.76%	69.24% +- 0.70%
Ft last 1	55.11% +- 0.83%	63.58% +- 0.75%	36.44% +- 0.78%	67.36% +- 0.71%
Ft last 2	55.65% +- 0.77%	63.44% +- 0.74%	36.24% +- 0.80%	67.31% +- 0.70%
Ft last 3	55.55% +- 0.83%	63.80% +- 0.76%	36.40% +- 0.73%	67.80% +- 0.66%
Mean centroid	55.25% +- 0.79%	62.64% +- 0.76%	40.62% +- 0.81%	65.40% +- 0.75%
Cosine Classifier	55.42% +- 0.77%	63.02% +- 0.73%	41.17% +- 0.85%	65.98% +- 0.70%
IMS	53.46% +- 0.85%	63.18% +- 0.77%	38.92% +- 0.79%	69.07% +- 0.75%



We have shown that fine-tuning methods outperform meta-learning methods for cross domain few-shot learning

How to choose which layers to fine-tune for a given dataset?



# Where to fine-tune in a deep network?

- Fine-tune just the last layer?
- Fine-tune the last  $K$  layers?
- Fine-tune all network parameters?
- Fine-tune a non-contiguous set of layers?
- How to make these choices for high capacity models with 10s, or 100s, or 1000s of layers?



# Where to fine-tune in a deep network?

- Fine-tune just the last layer?
- Fine-tune the last  $K$  layers?
- Fine-tune all network parameters?
- Fine-tune a non-contiguous set of layers?
- How to make these choices for high capacity models with 10s, or 100s, or 1000s of layers?

It depends on the dataset, pre-trained model, ...

**Fine-tuning is an art !**



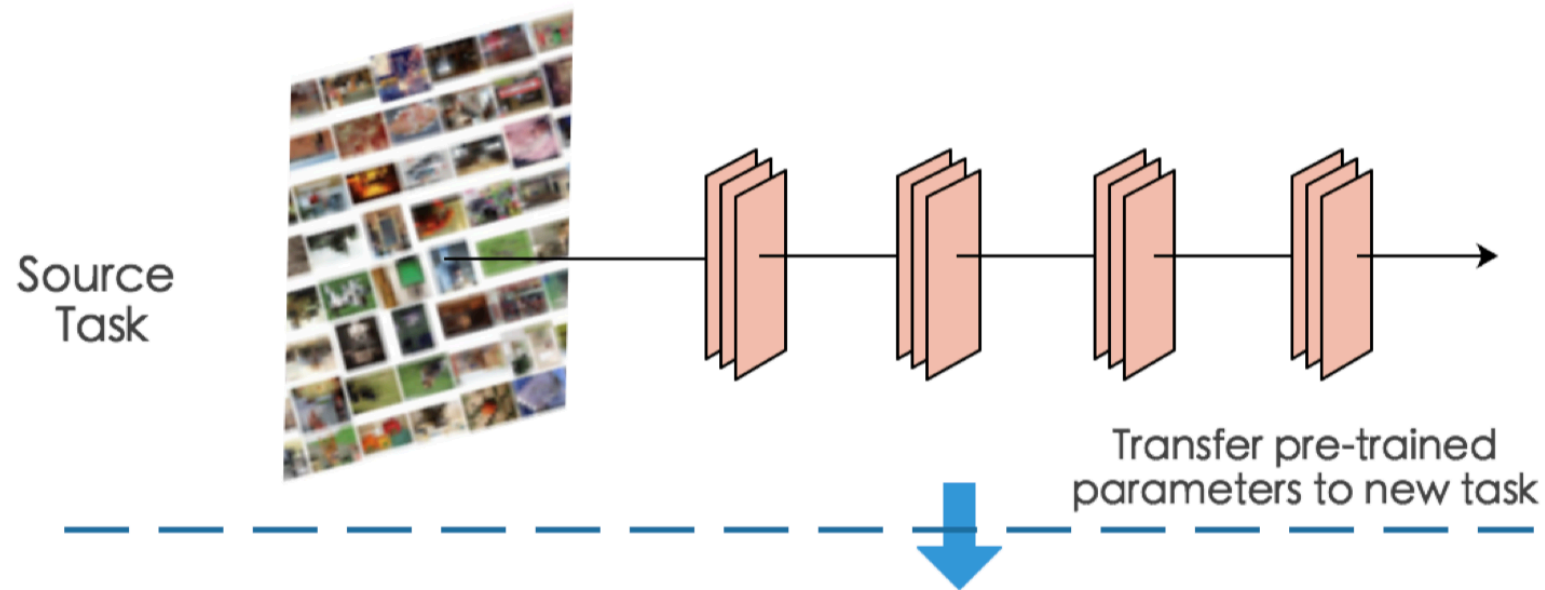


# SpotTune: Transfer Learning through Adaptive Fine-Tuning

Yunhui Guo, Honghui Shi, Abhishek Kumar, Kristen Grauman,  
Tajana Rosing, Rogerio Feris

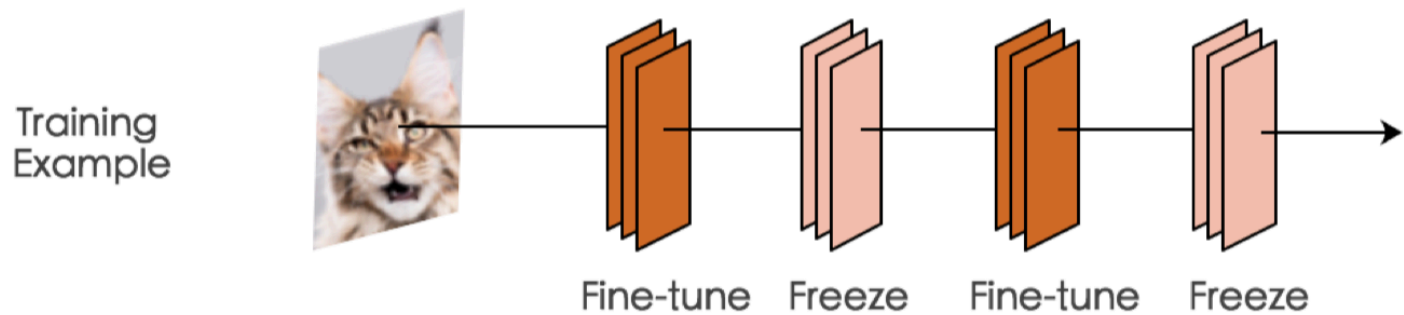
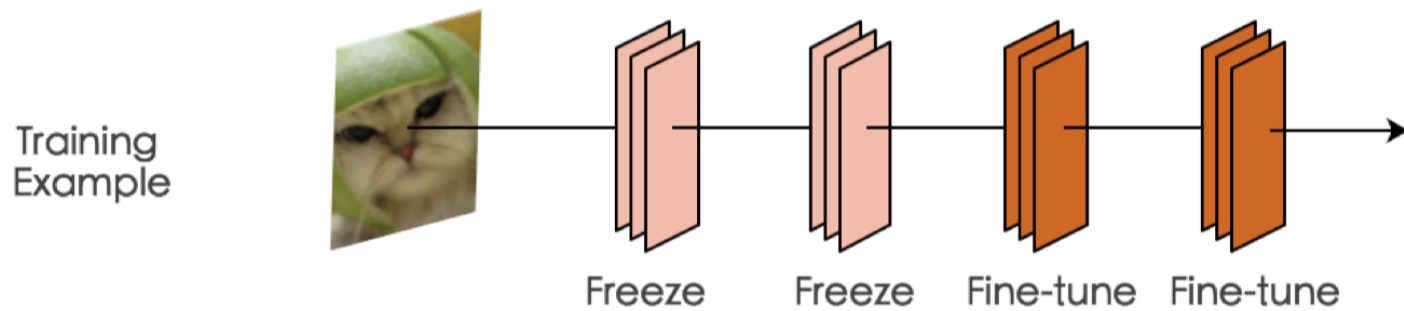
CVPR 2019



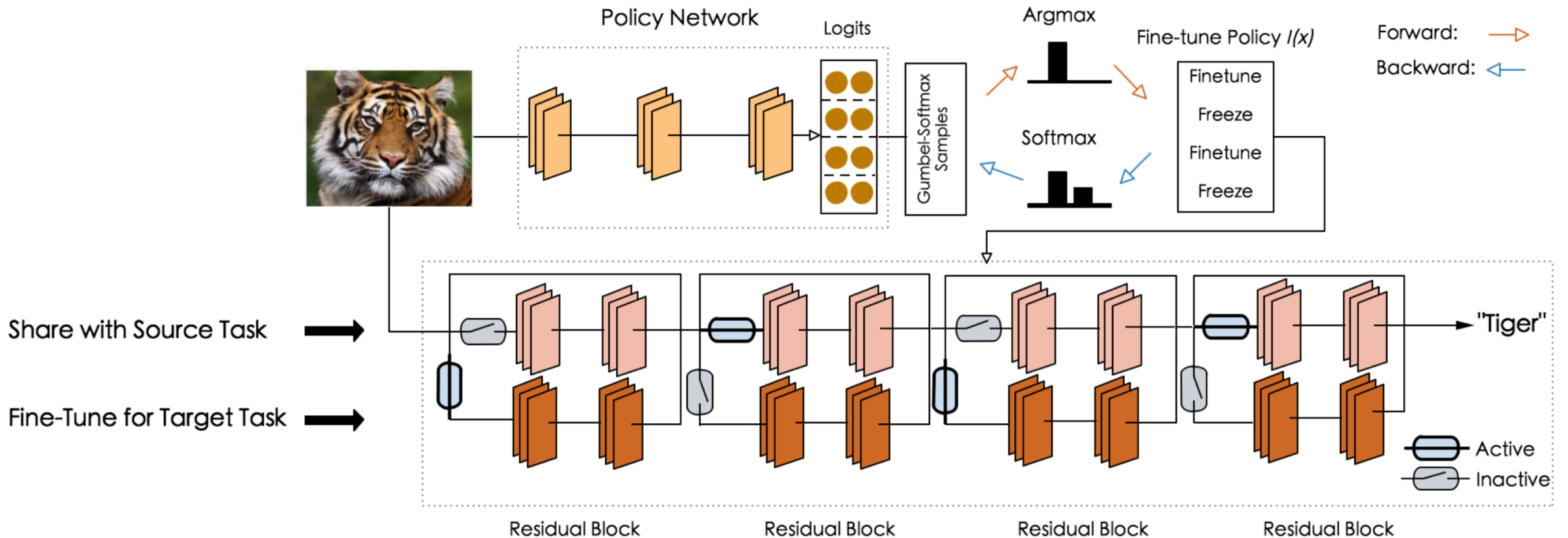


Target Task

Which layers to freeze and which layers to fine-tune?  
(per instance)

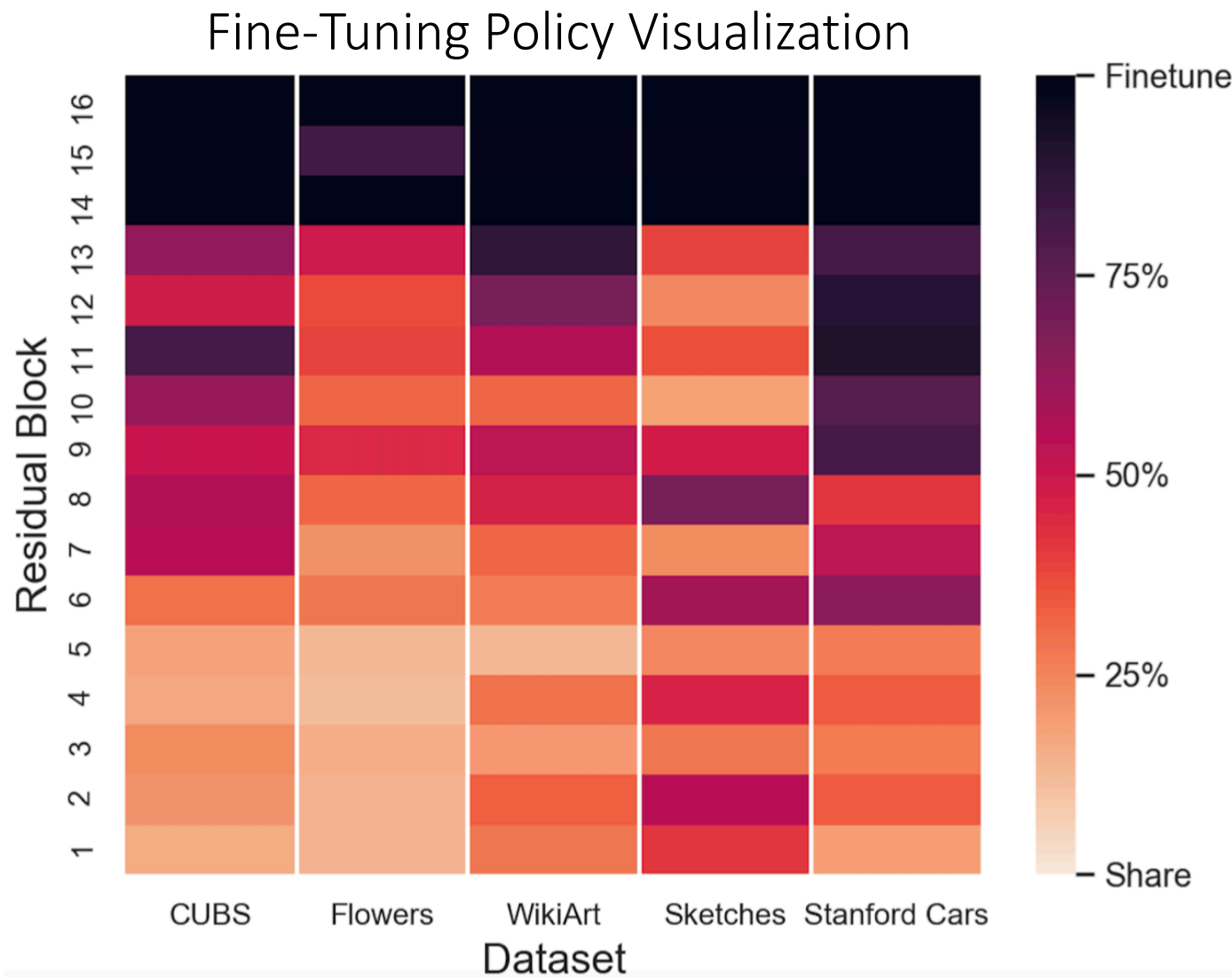


# SpotTune: Transfer Learning through Adaptive Fine-Tuning



\* General approach to any architecture (ResNet, VGG, ...)

# SpotTune: Transfer Learning through Adaptive Fine-Tuning



SpotTune automatically identifies the right fine-tuning policy for each dataset, for each training example.

# SpotTune: Transfer Learning through Adaptive Fine-Tuning

Model	CUBS	Stanford Cars	Flowers	WikiArt	Sketches
Feature Extractor	74.07%	70.81%	85.67%	61.60%	75.50%
Standard Fine-tuning	81.86%	89.74%	93.67%	75.60%	79.58%
Stochastic Fine-tuning	81.03%	88.94%	92.95%	73.06%	78.30%
Fine-tuning last-3	81.54%	88.21%	89.03%	72.68 %	77.72%
Fine-tuning last-2	80.34%	85.36%	91.81%	70.82%	78.37%
Fine-tuning last-1	78.68%	81.73%	89.99%	68.96%	77.20%
Random Policy	81.63 %	88.57%	93.44%	73.82%	78.30%
Fine-tuning ResNet-101	82.13%	90.32%	94.21%	<b>76.52%</b>	78.92%
$L^2$ -SP	83.69%	91.08%	95.21%	75.38%	79.60%
Progressive Neural Nets	83.08 %	91.59%	95.55%	75.41%	79.71%
SpotTune (running fine-tuned blocks)	82.36%	92.04%	93.49%	67.27%	78.88%
SpotTune (Global-k)	83.48%	90.51%	<b>96.60%</b>	75.63%	80.02%
SpotTune	<b>84.03 %</b>	<b>92.40 %</b>	96.34%	75.77%	<b>80.20%</b>





# SpotTune: Transfer Learning through Adaptive Fine-Tuning

	#par	ImNet	Airc.	C100	DPed	DTD	GTSR	Flwr	OGIt	SVHN	UCF	Score
Scratch	10x	59.87	57.10	75.73	91.20	37.77	96.55	56.30	88.74	96.63	43.27	1625
Scratch+ [37]	11x	59.67	59.59	76.08	92.45	39.63	96.90	56.66	88.74	96.78	44.17	1826
Feature Extractor	1x	59.67	23.31	63.11	80.33	55.53	68.18	73.69	58.79	43.54	26.80	544
Fine-tuning [38]	10x	60.32	61.87	82.12	92.82	55.53	99.42	81.41	89.12	96.55	51.20	3096
BN Adapt. [5]	1x	59.87	43.05	78.62	92.07	51.60	95.82	74.14	84.83	94.10	43.51	1353
LwF [26]	10x	59.87	61.15	82.23	92.34	58.83	97.57	83.05	88.08	96.10	50.04	2515
Series Res. adapt. [37]	2x	60.32	61.87	81.22	93.88	57.13	99.27	81.67	89.62	96.57	50.12	3159
Parallel Res. adapt. [38]	2x	60.32	64.21	81.92	94.73	58.83	99.38	84.68	89.21	96.54	50.94	3412
Res. adapt. (large) [37]	12x	67.00	67.69	84.69	94.28	59.41	97.43	84.86	89.92	96.59	52.39	3131
Res. adapt. decay [37]	2x	59.67	61.87	81.20	93.88	57.13	97.57	81.67	89.62	96.13	50.12	2621
Res. adapt. finetune all [37]	2x	59.23	63.73	81.31	93.30	57.02	97.47	83.43	89.82	96.17	50.28	2643
DAN [39]	2x	57.74	64.12	80.07	91.30	56.54	98.46	86.05	89.67	96.77	49.48	2851
PiggyBack [31]	1.28x	57.69	65.29	79.87	96.99	57.45	97.27	79.09	87.63	97.24	47.48	2838
SpotTune	11x	60.32	63.91	80.48	96.49	57.13	99.52	85.22	88.84	96.72	52.34	<b>3612</b>

SpotTune sets the new state of the art on the Visual Decathlon Challenge



# AdaShare: Learning What to Share for Efficient Multi-Task Learning

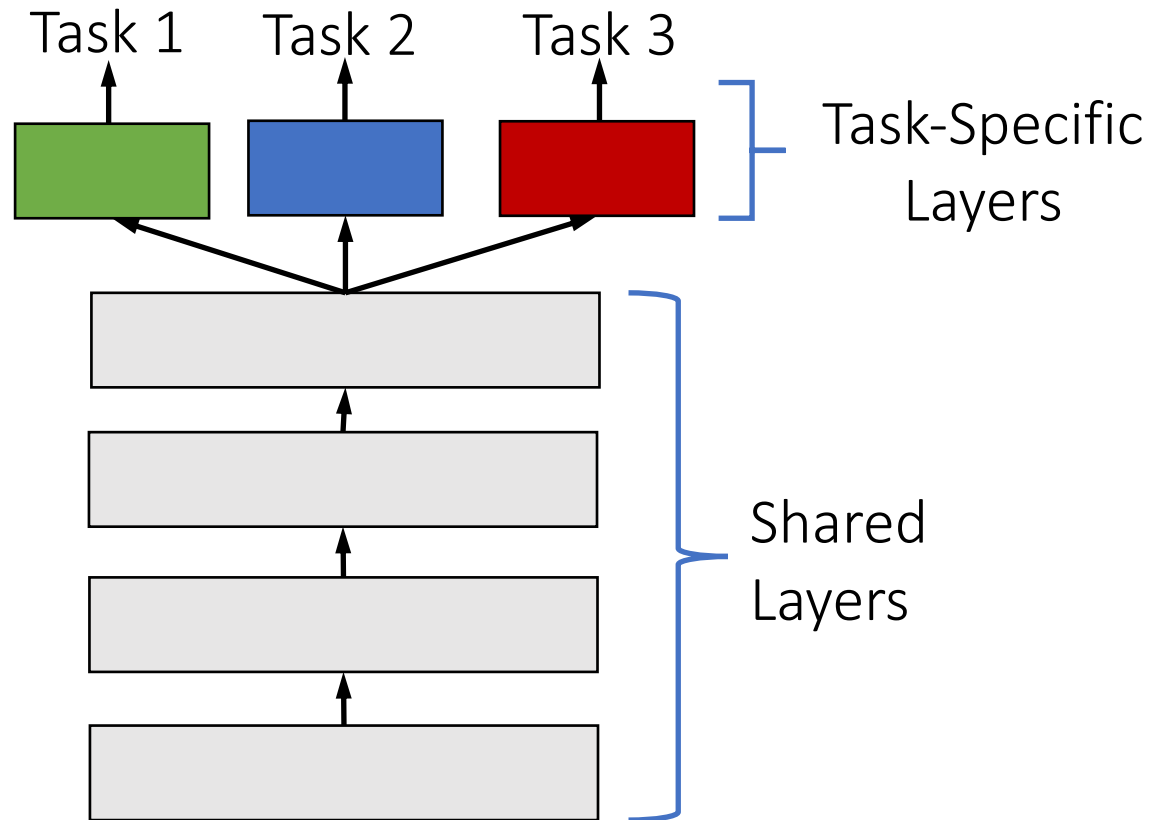
Ximeng Sun, Rameswar Panda, Rogerio Feris, Kate Saenko

NeurIPS 2020



# Hard Parameter Sharing

- Hand-designed architectures composed of base layers that are shared across tasks and specialized branches that learn task-specific features.



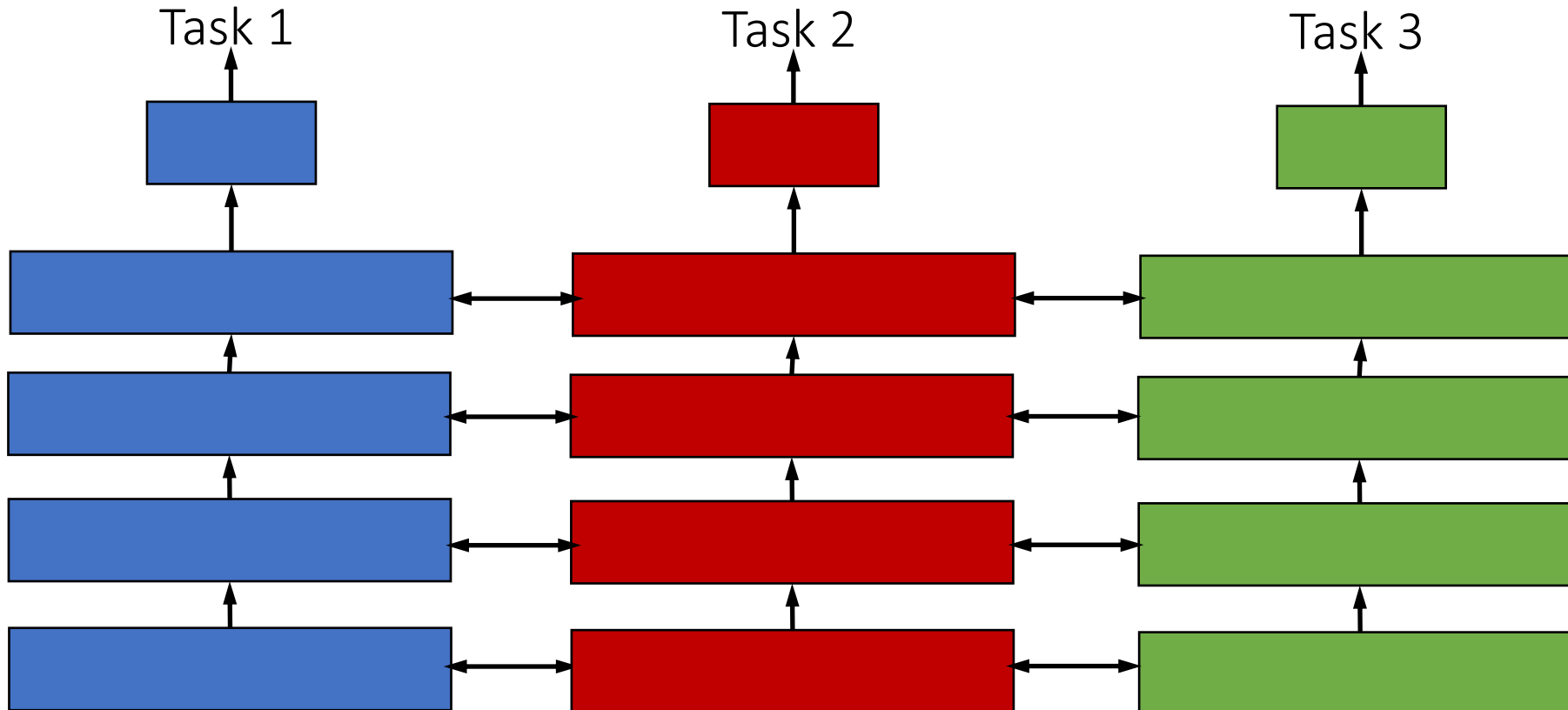
- Performance depends on “where to branch” in the network [Misra et al, 2016]
- The space of possible branching architectures is combinatorially large !



# Soft Parameter Sharing

- Network column for each task and a mechanism for feature sharing between columns.

Number of parameters grow linearly with the number of tasks !



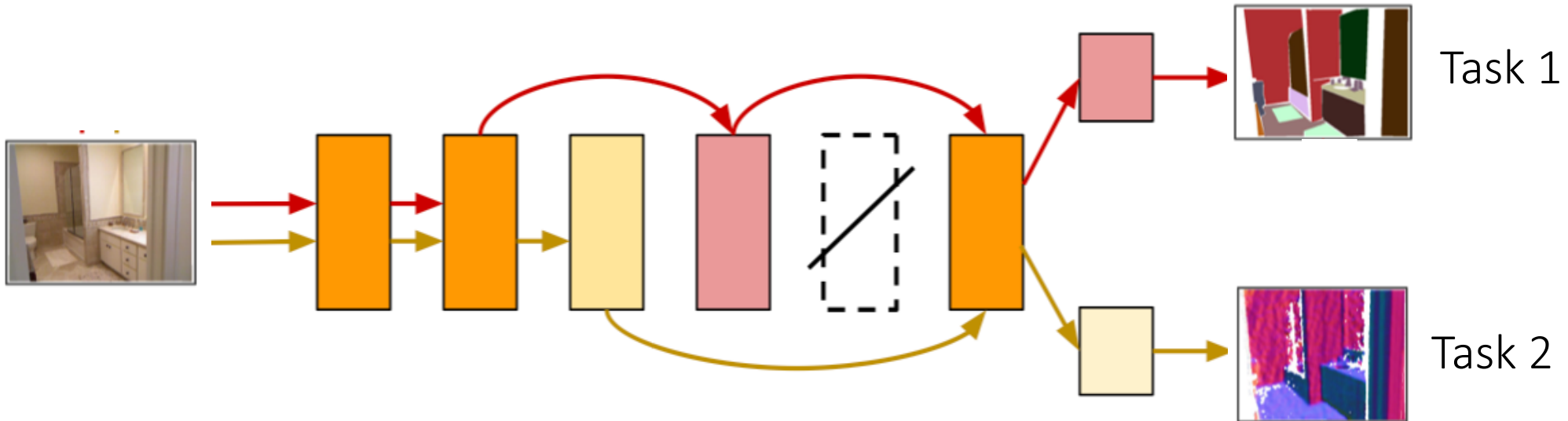
# Problem

*Can we determine which layers in the network should be shared across which tasks and which layers should be task-specific to achieve the best accuracy/memory footprint trade-off for scalable and efficient multi-task learning?*



# Proposed Approach: AdaShare

- Single network that supports separate execution paths for different tasks



Task 1-Specific



Task 2-Specific



Shared

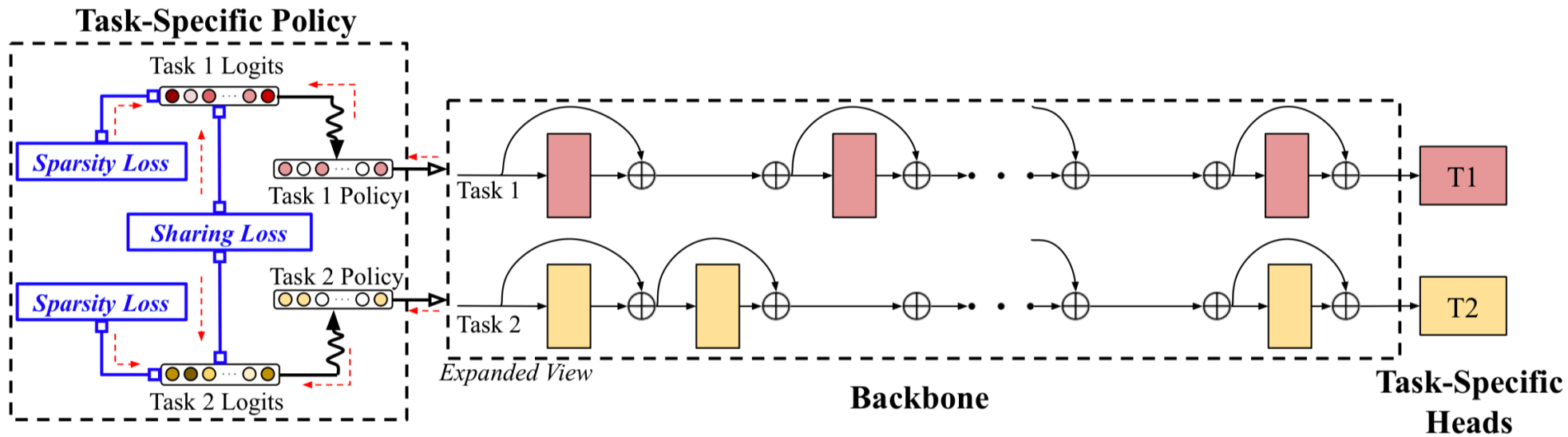


Skipped

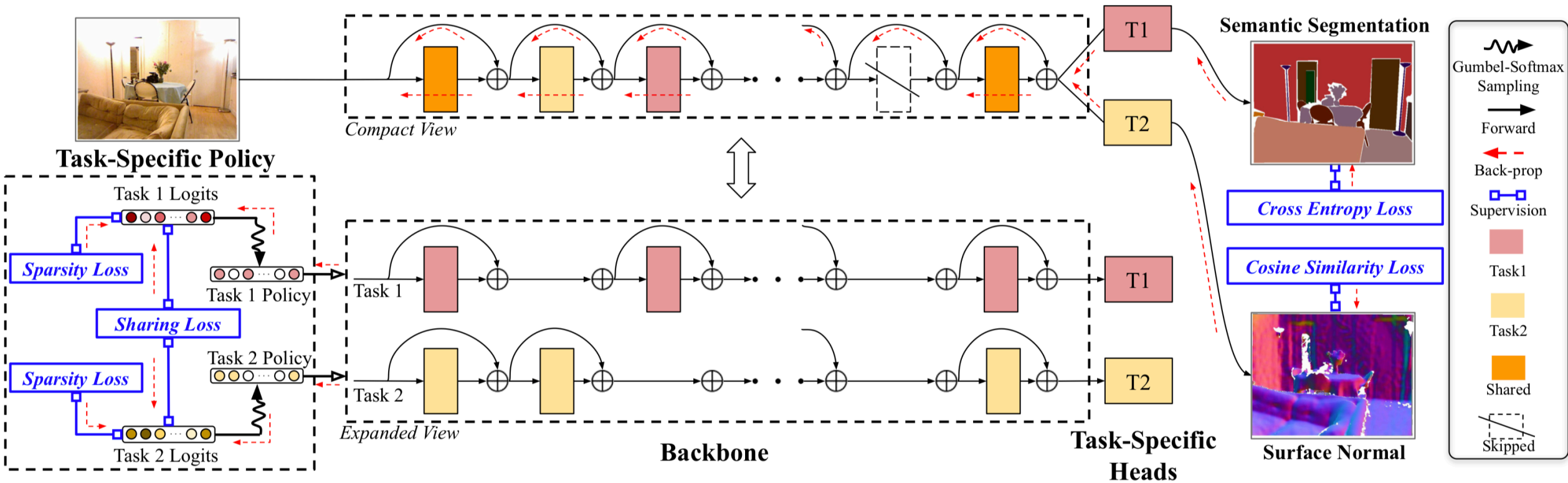




# AdaShare: Learning what to Share in Multi-Task Learning



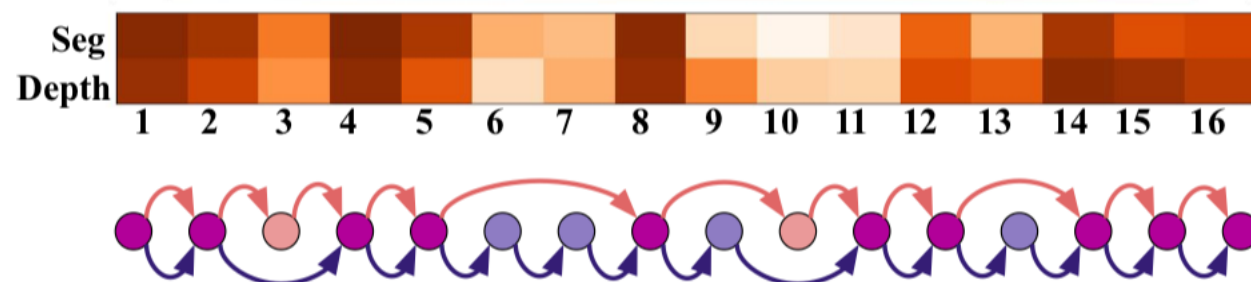
# AdaShare: Learning what to Share in Multi-Task Learning



# AdaShare: Experimental Results

- CityScapes [2 tasks]. *AdaShare* achieves the best performance on 5 out of 7 metrics using less than 1/2 parameters of most baselines.

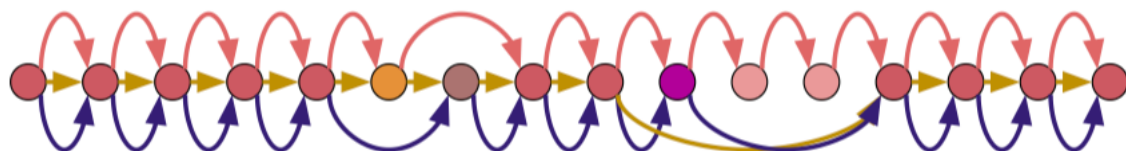
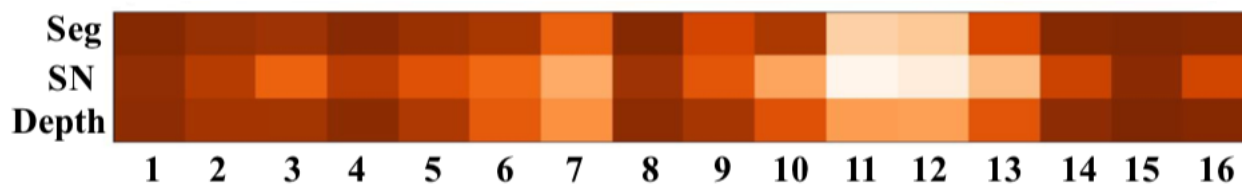
Model	# Params ↓	Semantic Seg.		Depth Prediction				
		mIoU ↑	Pixel Acc ↑	Error ↓		$\delta$ , within ↑		
				Abs	Rel	1.25	1.25 <sup>2</sup>	1.25 <sup>3</sup>
Single-Task	2	40.2	<u>74.7</u>	0.017	0.33	70.3	86.3	93.3
Multi-Task	<b>1</b>	37.7	73.8	0.018	0.34	72.4	88.3	94.2
Cross-Stitch	2	40.3	74.3	<b>0.015</b>	<b>0.30</b>	74.2	89.3	<b>94.9</b>
Sluice	2	39.8	74.2	<u>0.016</u>	<u>0.31</u>	73.0	88.8	94.6
NDDR-CNN	2.07	<b>41.5</b>	74.2	0.017	<u>0.31</u>	74.0	<u>89.3</u>	94.8
MTAN	2.41	<u>40.8</u>	74.3	<b>0.015</b>	0.32	<u>75.1</u>	89.3	94.6
<i>AdaShare</i>	<b>1</b>	<b>41.5</b>	<b>74.9</b>	<u>0.016</u>	0.33	<b>75.5</b>	<b>89.8</b>	<b>94.9</b>



# AdaShare: Experimental Results

- NYU v2 [3 tasks]. AdaShare achieves the best performance on 10 out of 12 metrics using less than 1/3 parameters of most baselines.

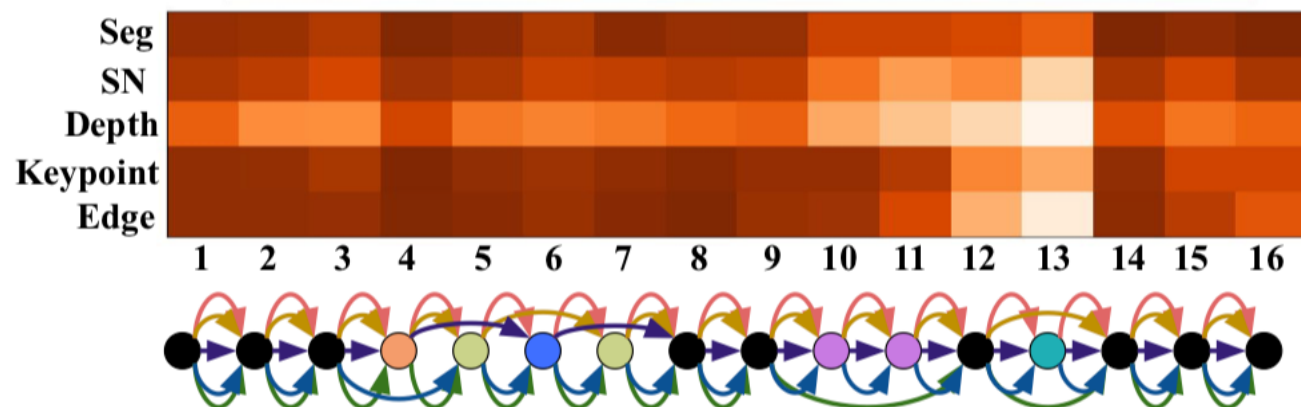
Model	# Params ↓	Semantic Seg.		Surface Normal Prediction					Depth Prediction				
		mIoU ↑	Pixel Acc ↑	Error ↓		$\theta$ , within ↑			Error ↓		$\delta$ , within ↑		
				Mean	Median	11.25°	22.5°	30°	Abs	Rel	1.25	1.25 <sup>2</sup>	1.25 <sup>3</sup>
Single-Task	3	<u>27.5</u>	<u>58.9</u>	17.5	15.2	34.9	<u>73.3</u>	85.7	0.62	0.25	57.9	85.8	95.7
Multi-Task	<b>1</b>	24.1	<u>57.2</u>	<b>16.6</b>	13.4	42.5	<u>73.2</u>	<u>84.6</u>	0.58	<u>0.23</u>	62.4	88.2	<u>96.5</u>
Cross-Stitch	3	25.4	57.6	17.2	14.0	41.4	70.5	<u>82.9</u>	0.58	<u>0.23</u>	61.4	<u>88.4</u>	<u>95.5</u>
Sluice	3	23.8	56.9	17.2	14.4	38.9	71.8	83.9	0.58	0.24	61.9	88.1	96.3
NDDR-CNN	3.15	21.6	53.9	<u>17.1</u>	14.5	37.4	<b>73.7</b>	<b>85.6</b>	0.66	0.26	55.7	83.7	94.8
MTAN	3.11	26.0	57.2	<b>16.6</b>	<u>13.0</u>	<u>43.7</u>	<u>73.3</u>	84.4	<u>0.57</u>	0.25	<u>62.7</u>	87.7	95.9
<i>AdaShare</i>	<b>1</b>	<b>30.2</b>	<b>62.4</b>	<b>16.6</b>	<b>12.9</b>	<b>45.0</b>	71.7	83.0	<b>0.55</b>	<b>0.20</b>	<b>64.5</b>	<b>90.5</b>	<b>97.8</b>



# AdaShare: Experimental Results

- **Tiny-Taskonomy [5 Tasks]**. AdaShare outperforms the baselines on 3 out of 5 tasks using less than 1/5 parameters of most baselines.

Models	# Params ↓	Seg ↓	SN ↑	Depth ↓	Keypoint ↓	Edge ↓
Single-Task	5	0.575	<b>0.707</b>	<b>0.022</b>	0.197	0.212
Multi-Task	<b>1</b>	0.587	0.702	0.024	0.194	0.201
Cross-Stitch	5	<u>0.560</u>	0.684	<b>0.022</b>	0.202	0.219
Sluice	5	0.610	0.702	<u>0.023</u>	<b>0.192</b>	<u>0.198</u>
NDDR-CNN	5.41	<b>0.539</b>	<u>0.705</u>	0.024	0.194	0.206
MTAN	4.51	0.637	0.702	<u>0.023</u>	<u>0.193</u>	0.203
<i>AdaShare</i>	<b>1</b>	0.566	<b>0.707</b>	0.025	<b>0.192</b>	<b>0.193</b>





# Visual Learning Beyond Natural Images: Summary

- Naïve fine-tuning outperforms current meta-learning approaches for cross-domain (beyond natural images) few-shot learning
- The optimal set of layers to fine-tune is dependent on the dataset. SpotTune automatically decides which layers of a model should be shared with the pre-trained model and which layers should be fine-tuned
- Deciding what features should be shared is also crucial for joint multi-task learning. AdaShare selects specific computational paths for each task to maximize accuracy and efficiency.



# Thank you !

- Y. Guo, N. Codella, L. Karlinsky, J. Smith, T. Rosing and R. S. Feris. "A new Benchmark for Evaluation of Cross-domain Few-shot Learning." ECCV 2020 (\* equal contribution) [\[code available\]](#)
- Y. Guo, H. Shi, A. Kumar, K. Grauman, T. Rosing and R. S. Feris. "SpotTune: Transfer Learning Through Adaptive Fine-Tuning" CVPR 2019 [\[code available\]](#)
- X. Sun, R. Panda and R. S. Feris. "AdaShare: Learning What to Share for Efficient Deep Multi-Task Learning". NeurIPS 2020

