



Optimization and Generalization of Neural Networks at the Edge

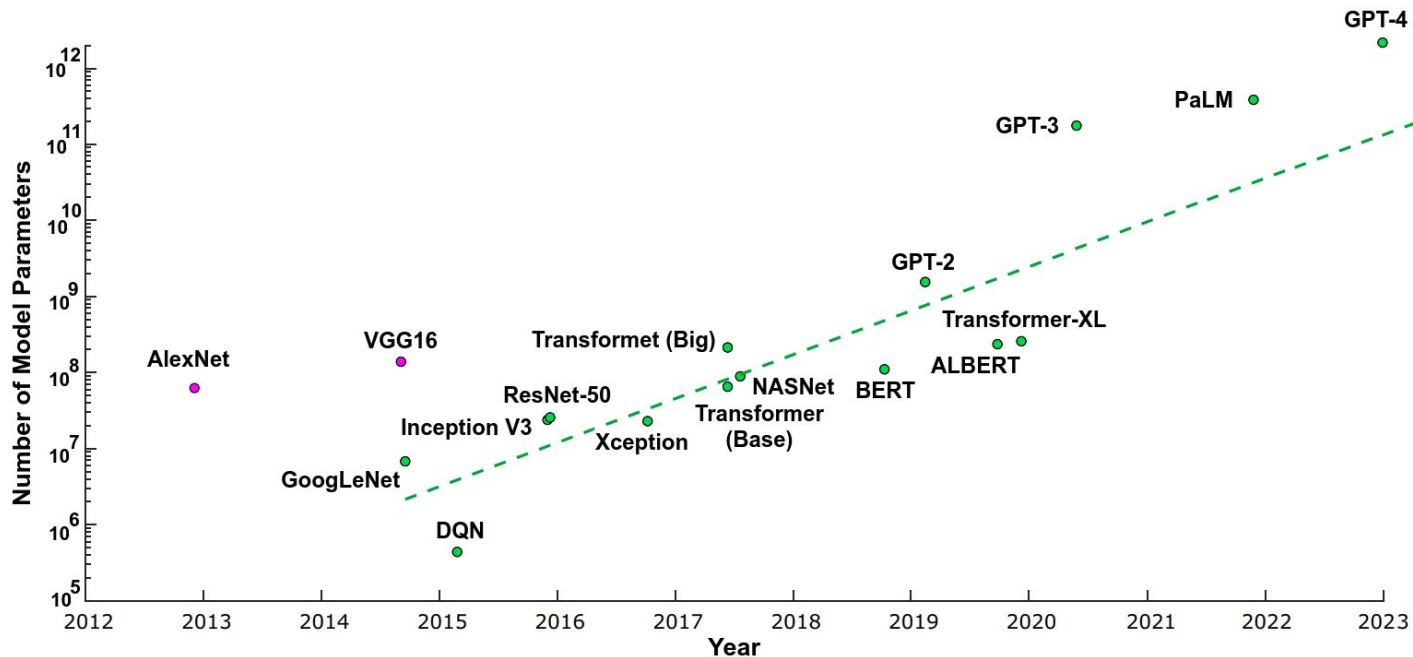
Rahim Entezari

July 17th, 2023

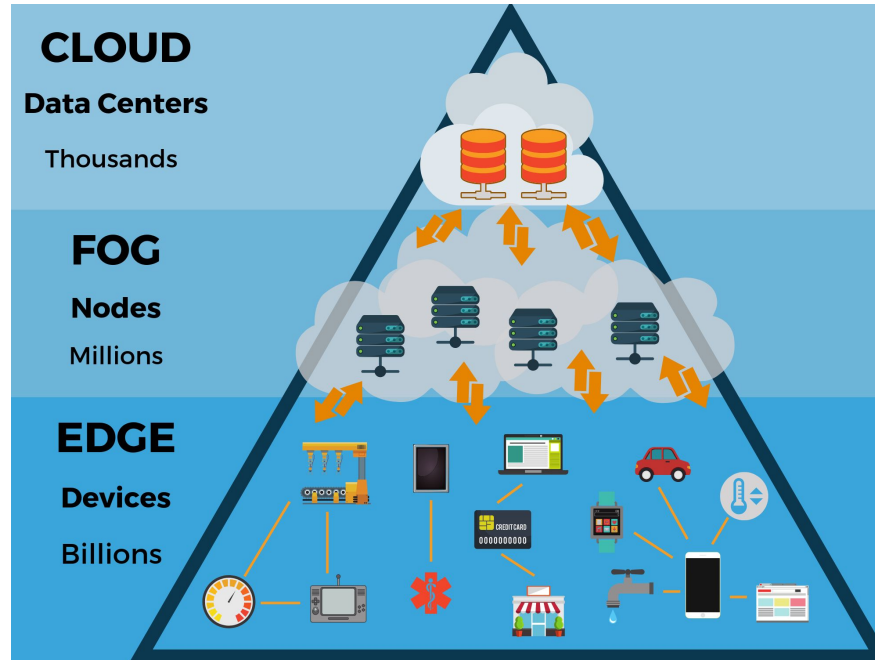


COMPLEXITY
SCIENCE
HUB
VIENNA

Scaling trend

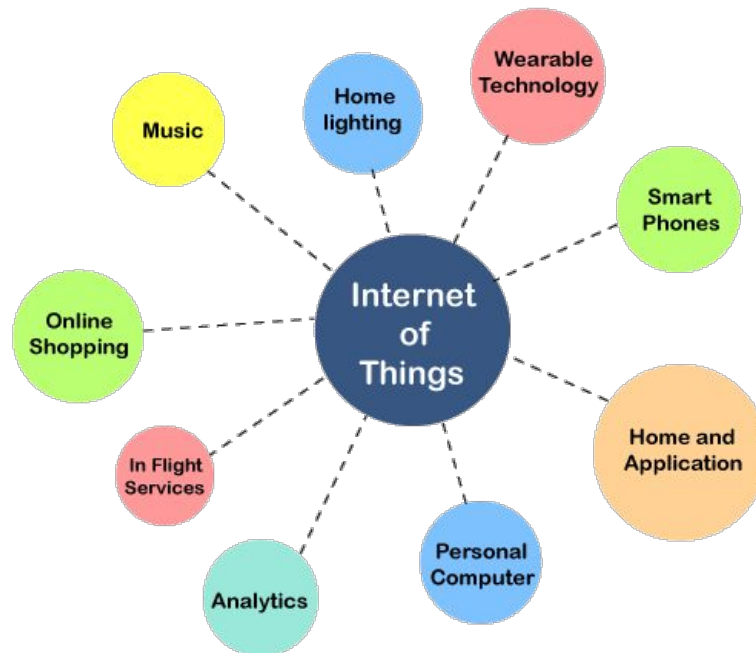


Neural networks at the edge



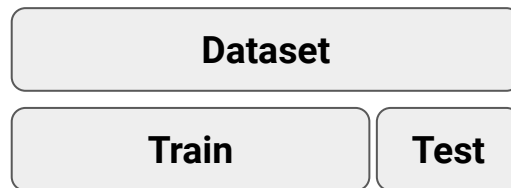
Neural networks at the edge

Pervasive but limited resources → make AI possible on the edge



Generalization

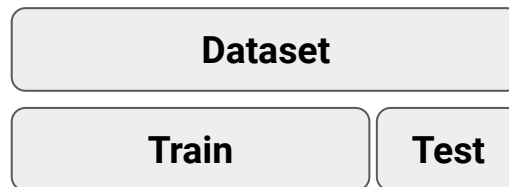
In-Distribution vs. Out-Of-Distribution generalization (ID vs. OOD)



ID generalization

Generalization

In-Distribution vs. Out-Of-Distribution generalization (ID vs. OOD)

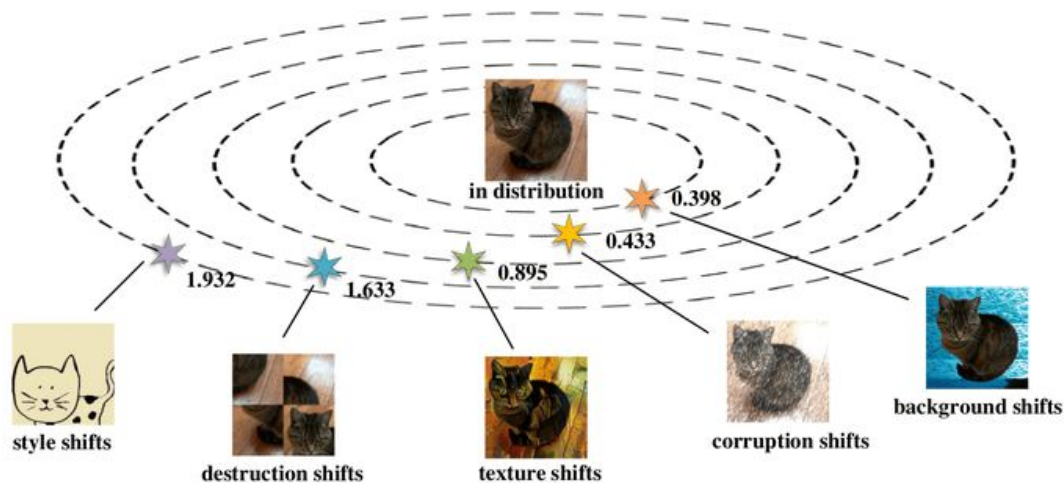


ID generalization

- Regularization
- Dropout
- Early stopping

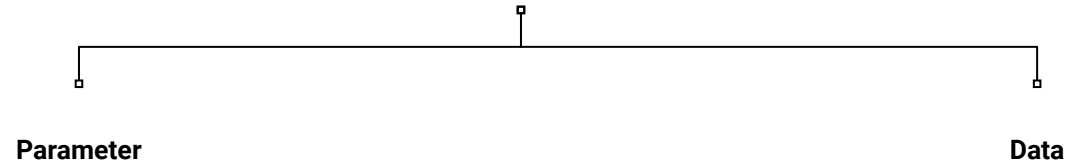
Generalization

In-Distribution vs. Out-Of-Distribution generalization (ID vs. OOD)



- Transfer learning
- Domain adaptation

Generalization of Neural Networks

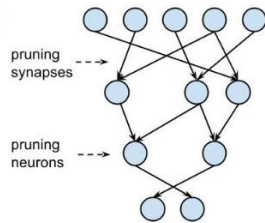


Generalization of Neural Networks

Parameter

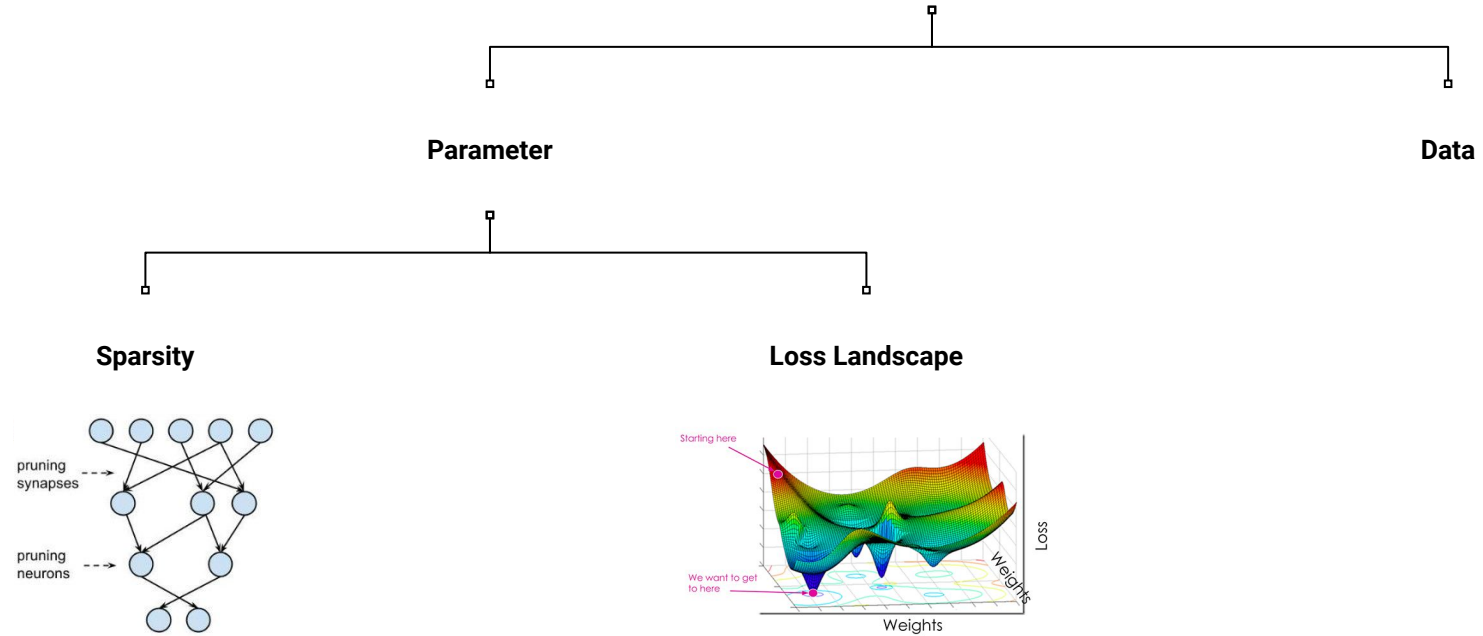
Data

Sparsity

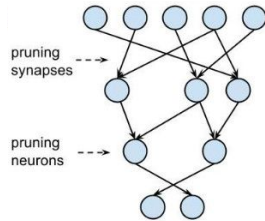


Make neural networks work at the edge
To improve generalization

Generalization of Neural Networks

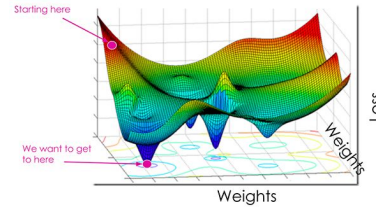


Sparsity



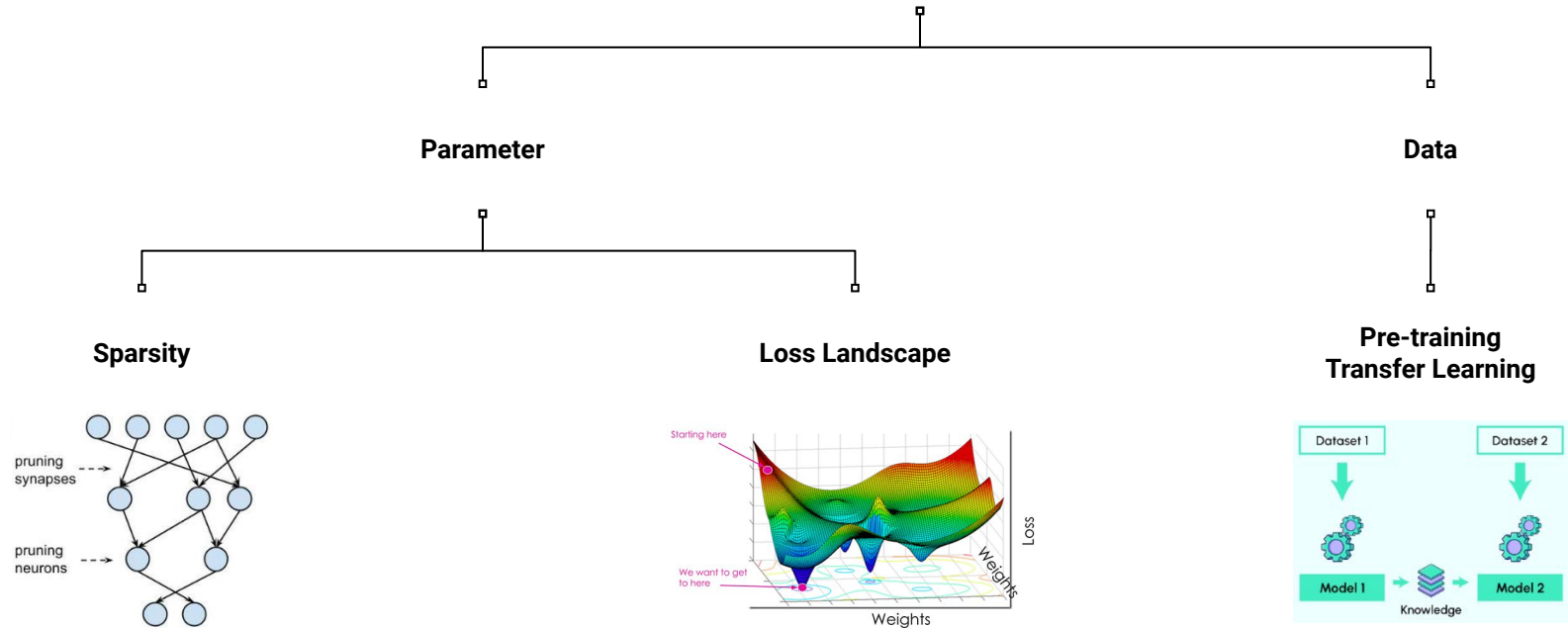
Make neural networks work at the edge
To improve generalization

Loss Landscape



To understand/probe trained networks
ensembles to make neural networks work at the edge

Generalization of Neural Networks

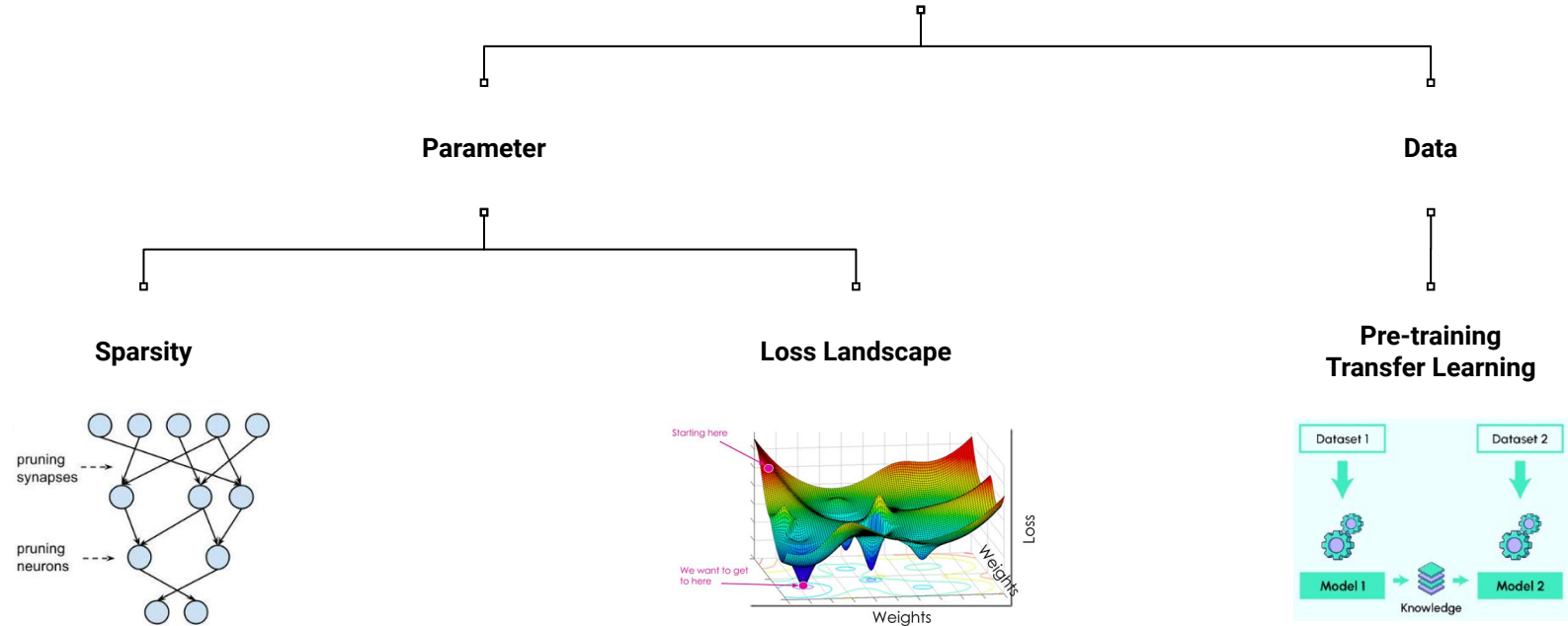


Make neural networks work at the edge
To improve generalization

To understand/probe trained networks
ensembles to make neural networks work at the edge

Dynamic environment → dynamic data
Reliable AI applications

Generalization of Neural Networks



Class-dependent pruning of deep neural networks

Understanding the effect of sparsity on neural networks robustness

Studying the impact of magnitude pruning on contrastive learning methods

The Role of Permutation Invariance in Linear Mode Connectivity of Neural Networks

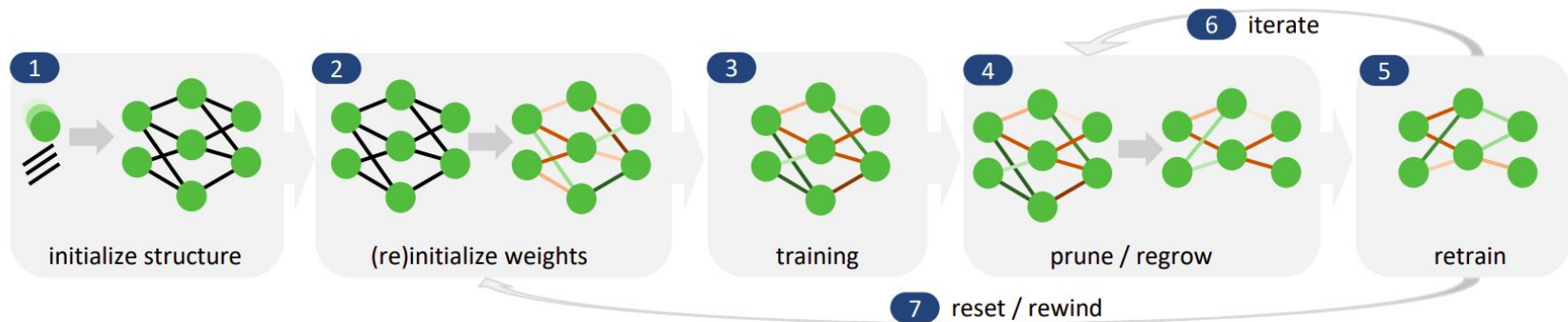
REPAIR: Linear Mode Connectivity of Deep Neural Networks via Permutation Invariance and Renormalization

The Role of Pretraining Data in Transfer Learning

Part 1: Sparsity

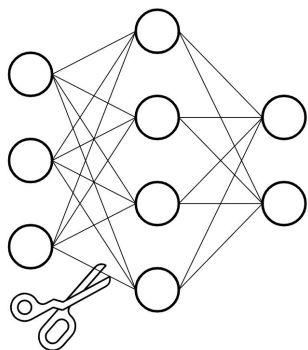
Sparsity

“With all things being equal, the simplest explanation tends to be the right one” (William of Ockham, ~1300)

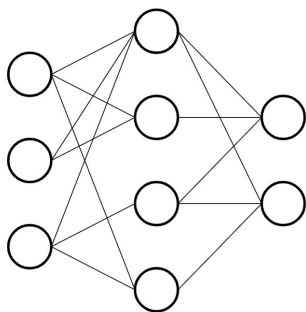


Relation between sparsity and generalization

Does sparsity help/hurt generalization?

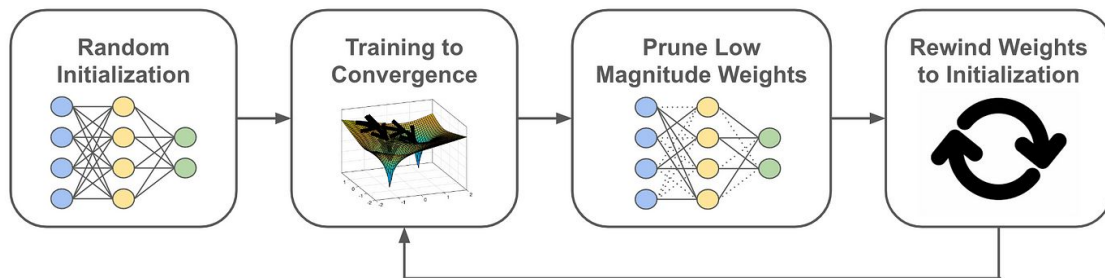


Before pruning



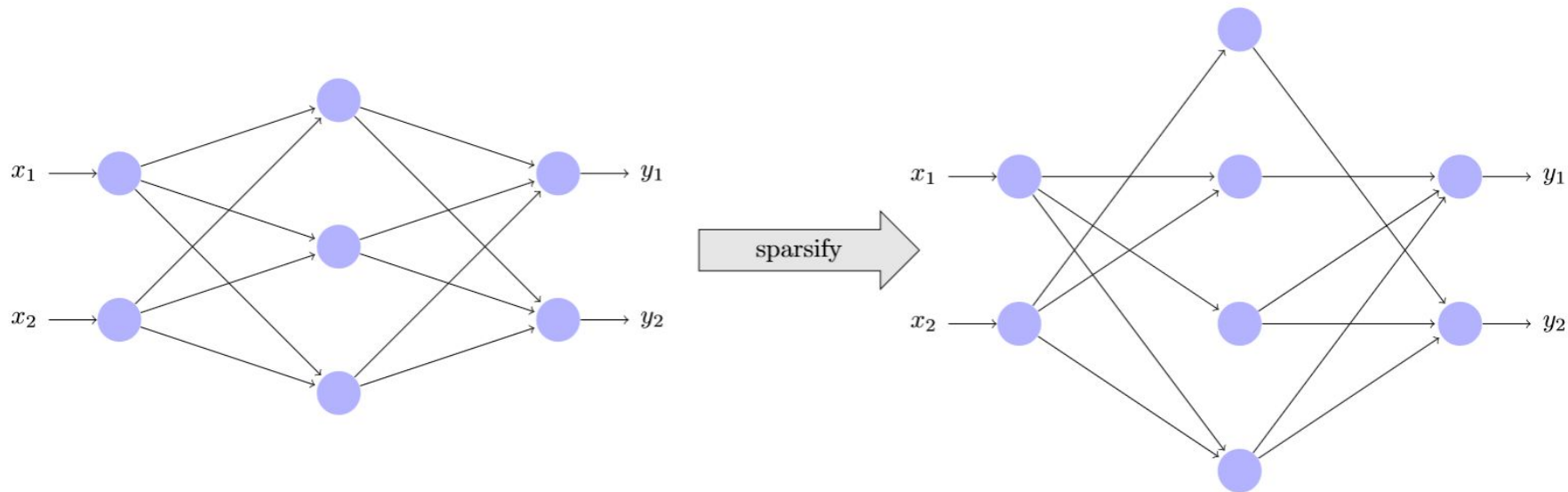
After pruning

Magnitude pruning



Lottery Ticket

Effective model capacity



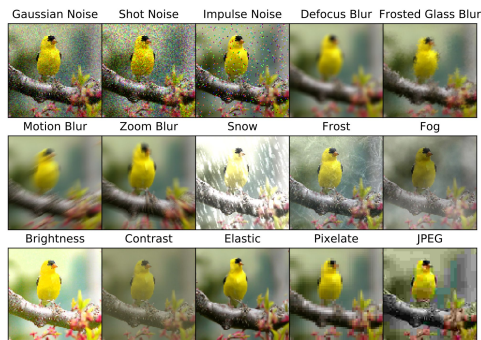
Sparsity and Generalization

1. Data corruption

2. Weight perturbation

Sparsity and Generalization

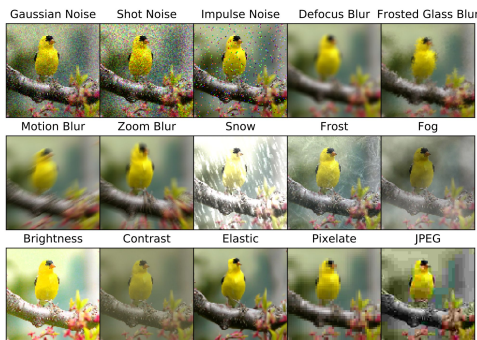
1. Data corruption



- Performance on corrupted datasets
 - MNIST-C
 - CIFAR10-C
 - CIFAR100-C

Sparsity and Generalization

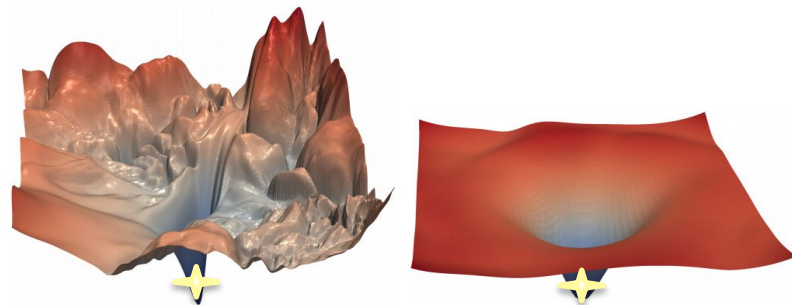
1. Data corruption



- Performance on corrupted datasets
 - MNIST-C
 - CIFAR10-C
 - CIFAR100-C

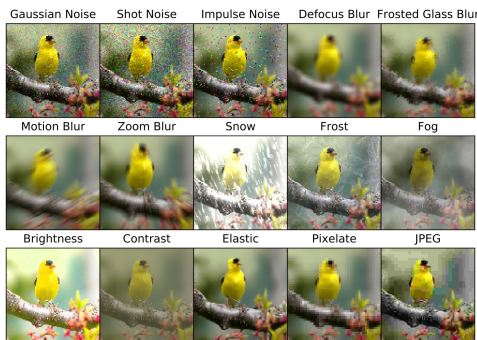
2. Weight perturbation

- Add Gaussian noise to each weight
 - $z_i \sim N(\mu, \omega_i^2 \sigma_i^2)$
- Flatness of achieved minima



Sparsity and Generalization

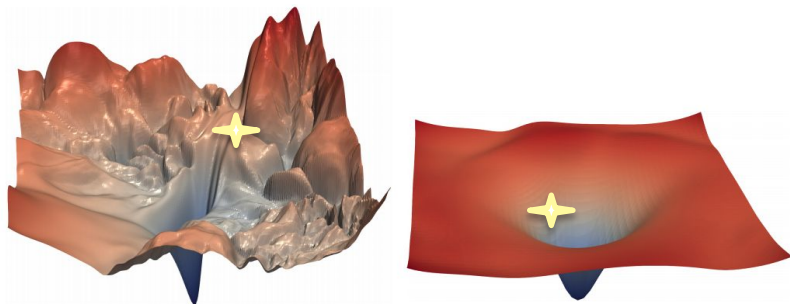
1. Data corruption



- Performance on corrupted datasets
 - MNIST-C
 - CIFAR10-C
 - CIFAR100-C

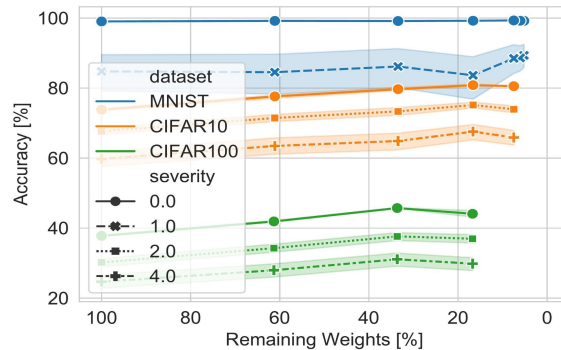
2. Weight perturbation

- Add Gaussian noise to each weight
 - $z_i \sim N(\mu, \omega_i^2 \sigma_i^2)$
- Flatness of achieved minima

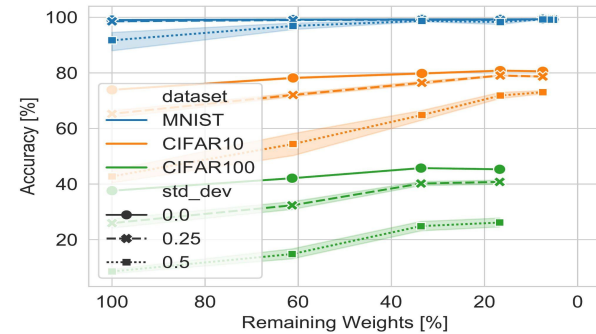


Sparsity and Generalization

1. Data corruption



2. Weight perturbation

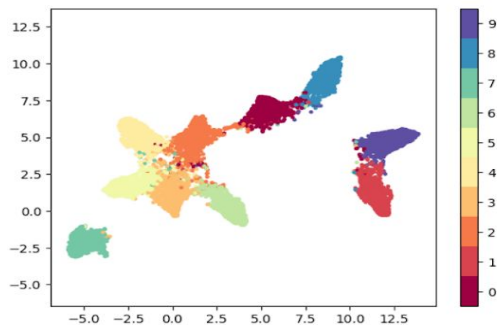


contrary to common belief, sparsity indeed does not hurt network generalization

What is the effect of sparsity on learned representations?

Learned representations: UMAP

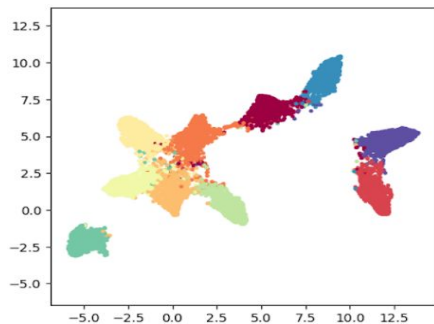
Supervised



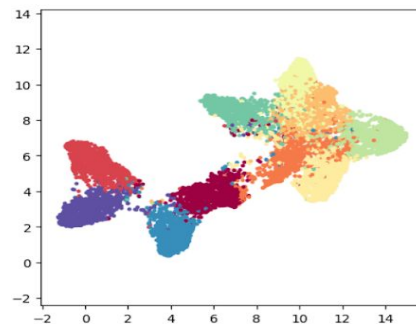
Dense

Learned representations: UMAP

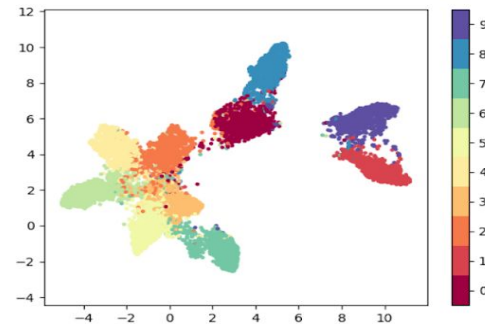
Supervised



Dense



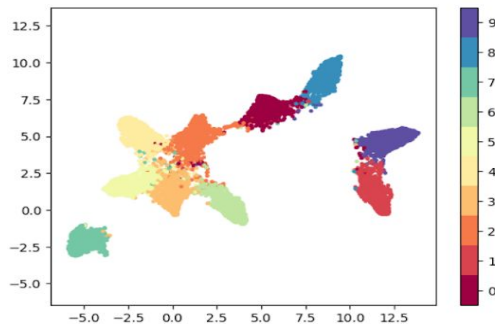
GMP 90%



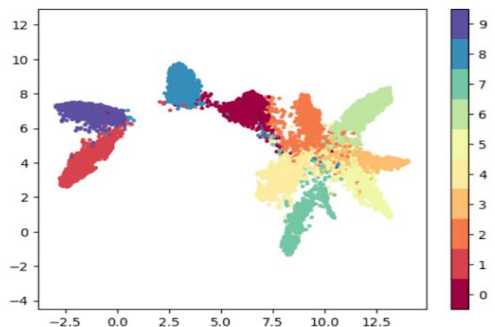
One-shot 90%

UMAP: what if we change the training algorithm?

Supervised



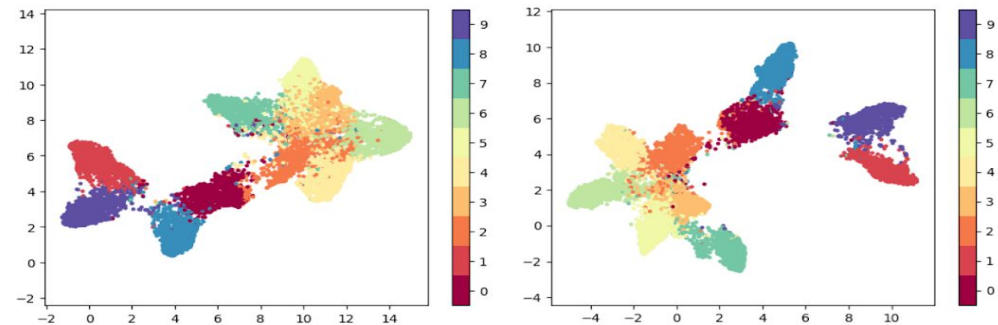
**Supervised
Contrastive**



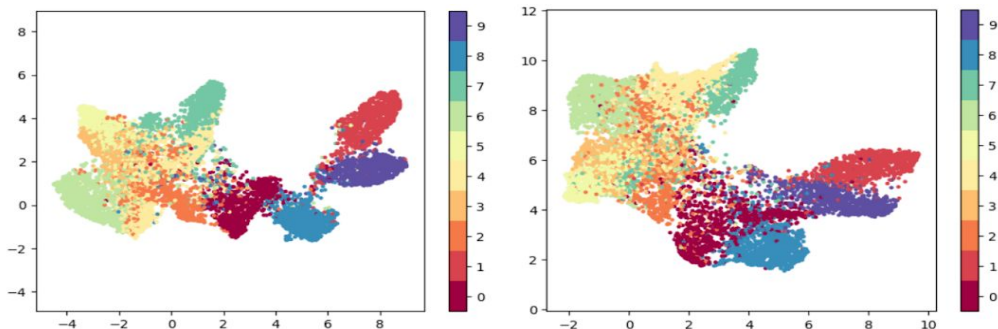
Dense

UMAP: supervised vs. semi-supervised

Supervised



**Supervised
Contrastive**



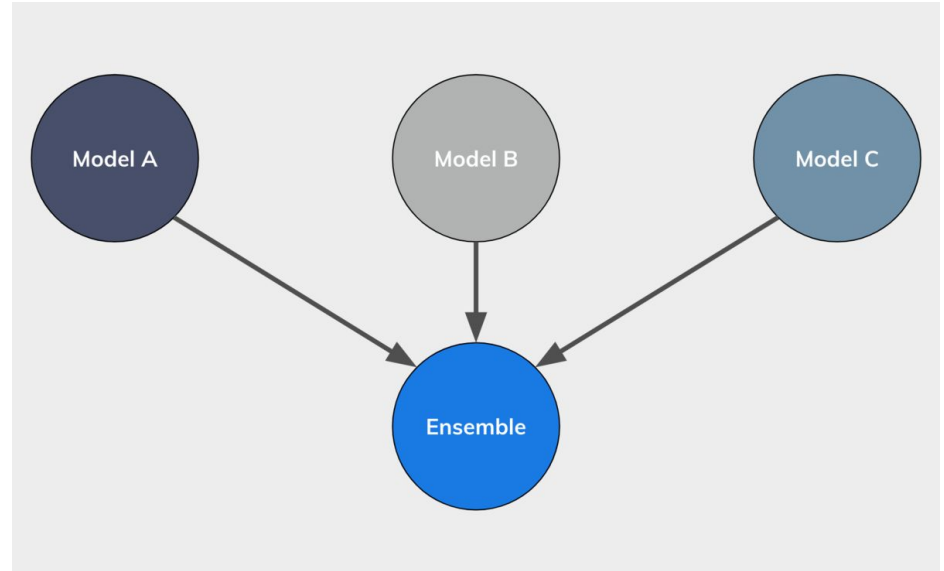
GMP 90%

One-shot 90%

Part 2: Loss Landscape

Motivation

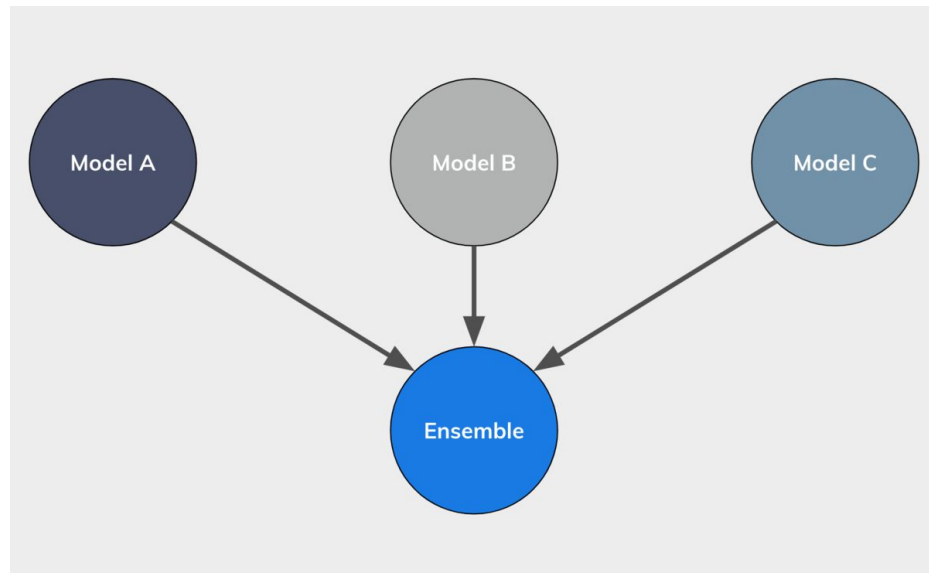
Ensembling helps generalization



Motivation

Form an ensemble model

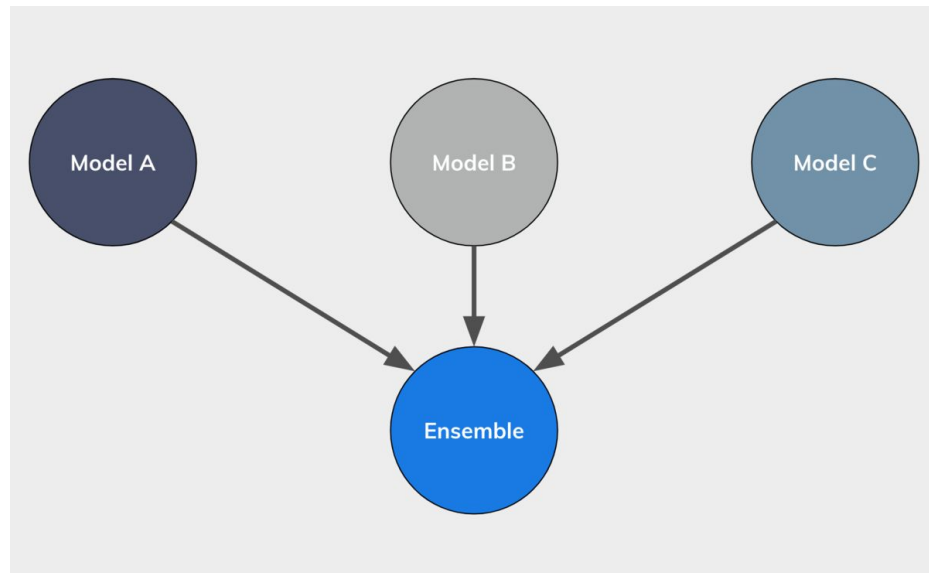
1. In output space



Motivation

Form an ensemble model

1. In output space
2. **In weight space (Embedded ML)**



Weight Averaging

Ensemble by weight averaging

Requirements:

1. Solutions should be functionally diverse

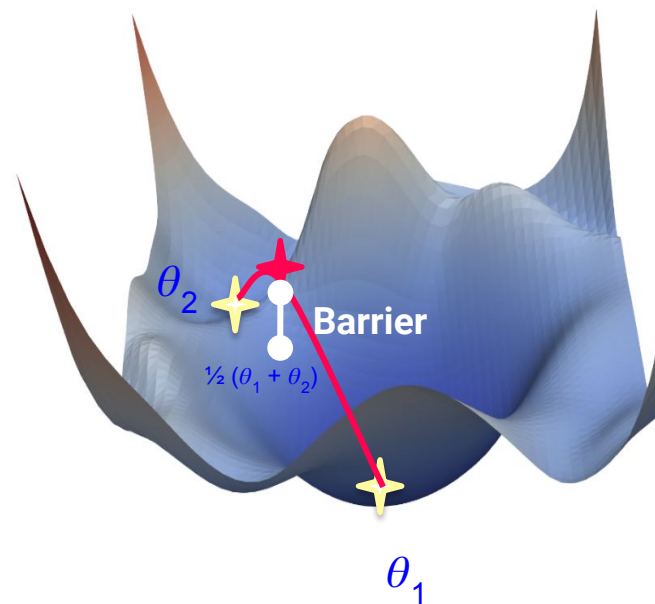


Weight Averaging

Ensemble by weight averaging

Requirements:

1. Solutions should be functionally diverse
2. Solutions should reside in one basin

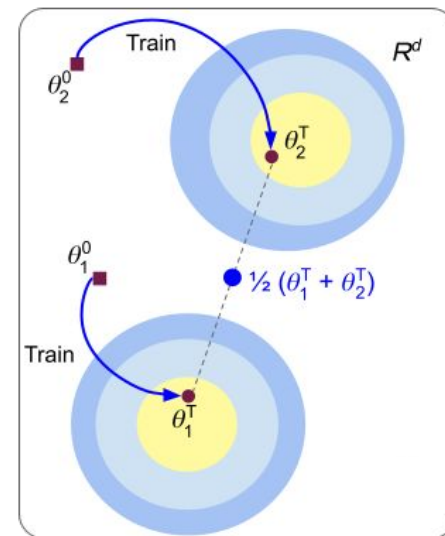


Related works

Functionally different solutions



Weight space averaging

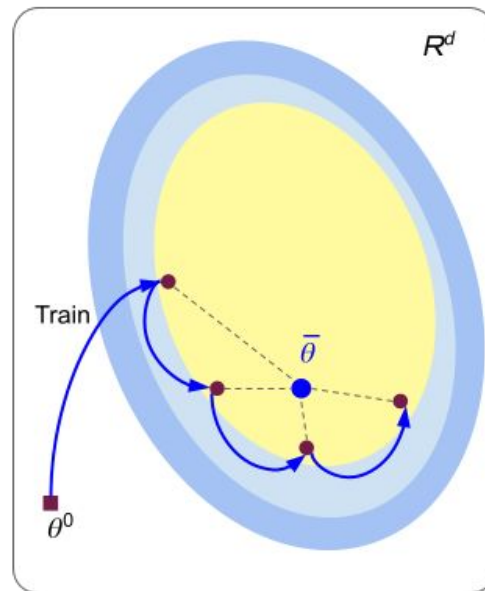


Related works

Functionally different solutions



Weight space averaging



Is there any way to make different solutions in one basin?

Functionally different solutions

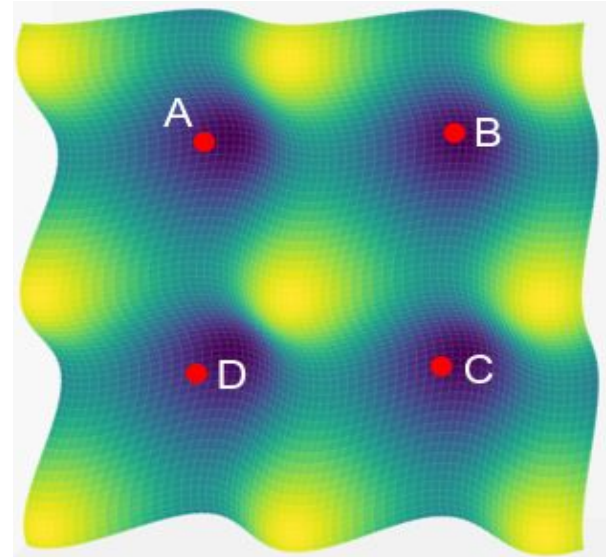


Weight space averaging



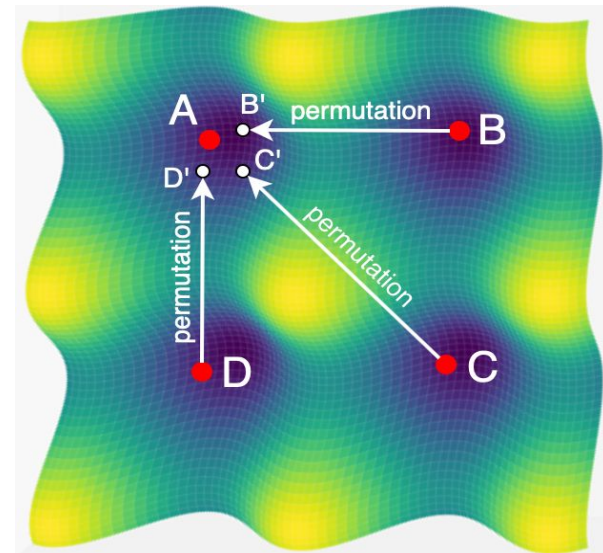
Conjecture

A, B, C, and D are minimas in different basins with barriers between pairs.



Conjecture

Taking permutations into account, there is likely no barrier in the linear interpolation between SGD solutions.



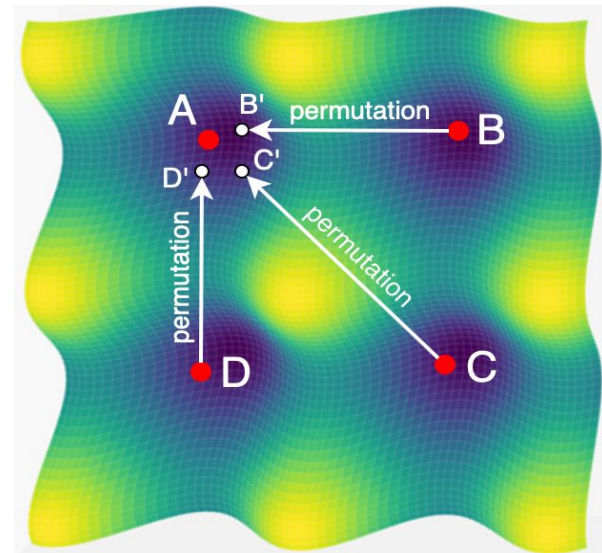
Conjecture

Taking permutations into account, there is likely no barrier in the linear interpolation between SGD solutions.

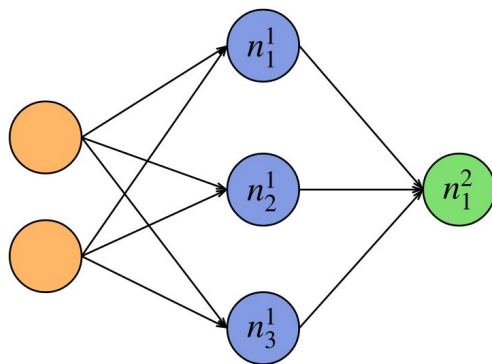
Functionally different solutions



Weight space averaging

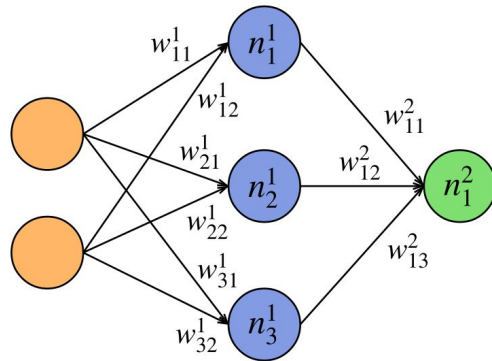


Permutations in Neural Networks

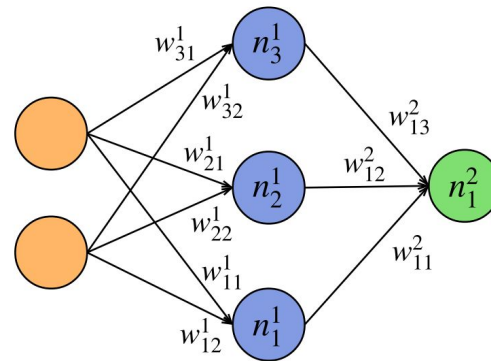


$3! = 6$ permutations

Permutations does not change the function!



(a) Neural network f_1



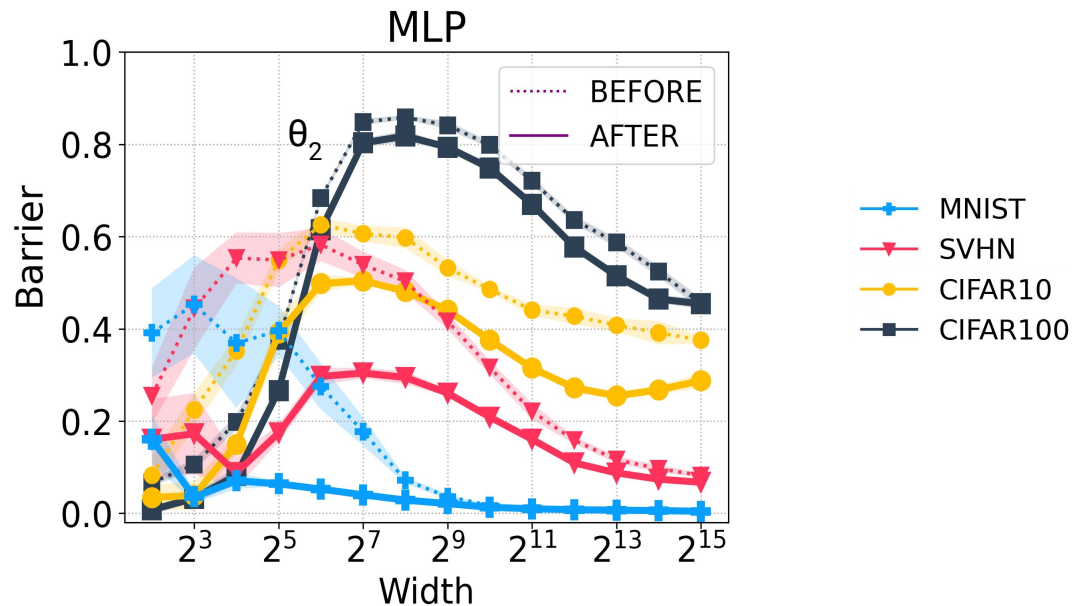
(b) Neural network f_2

How to find the appropriate permutation?

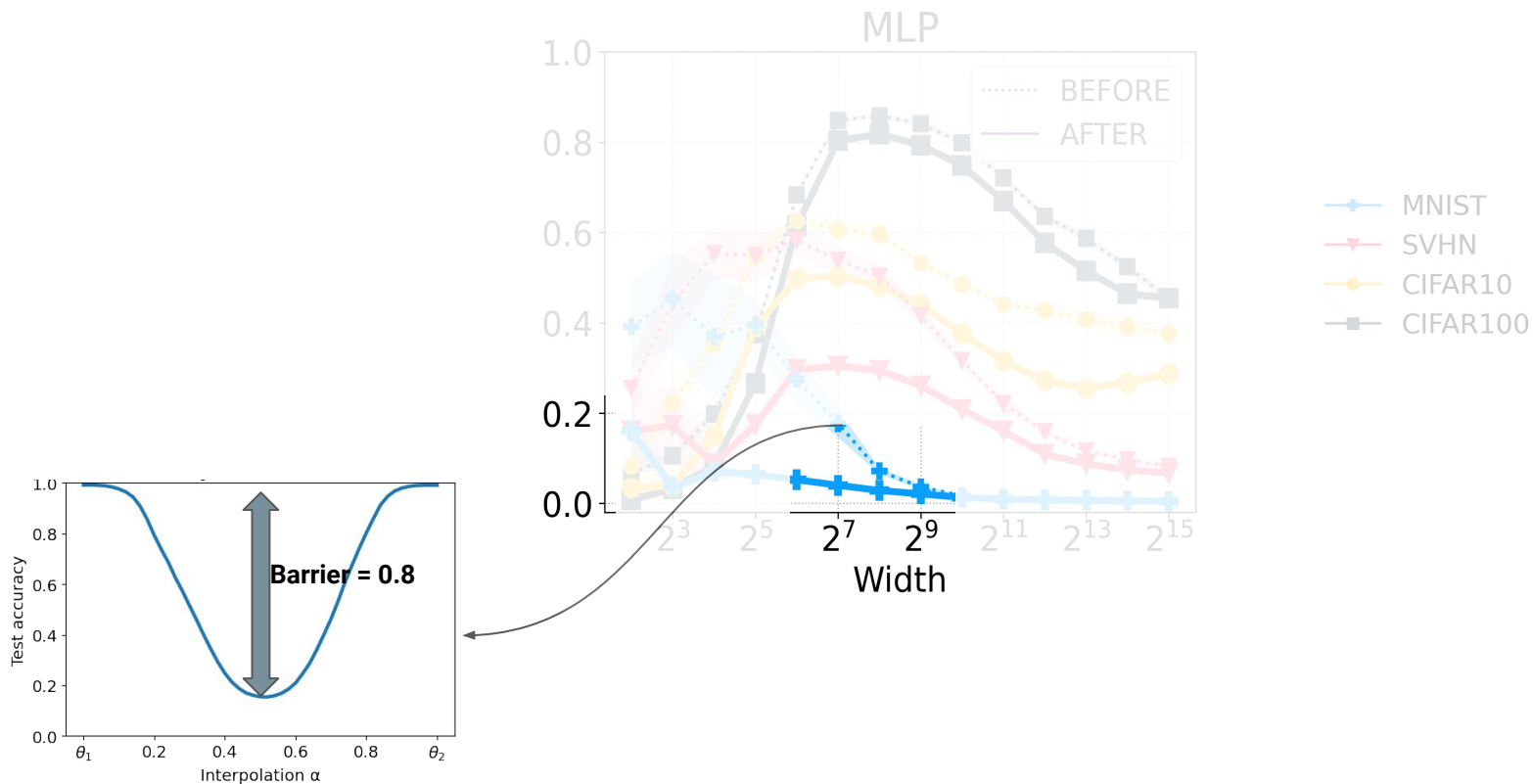
Permutation by brute force

- ResNet-50 $\rightarrow 10^{55109}$
- For comparison, the number of atoms in universe is about 10^{82}

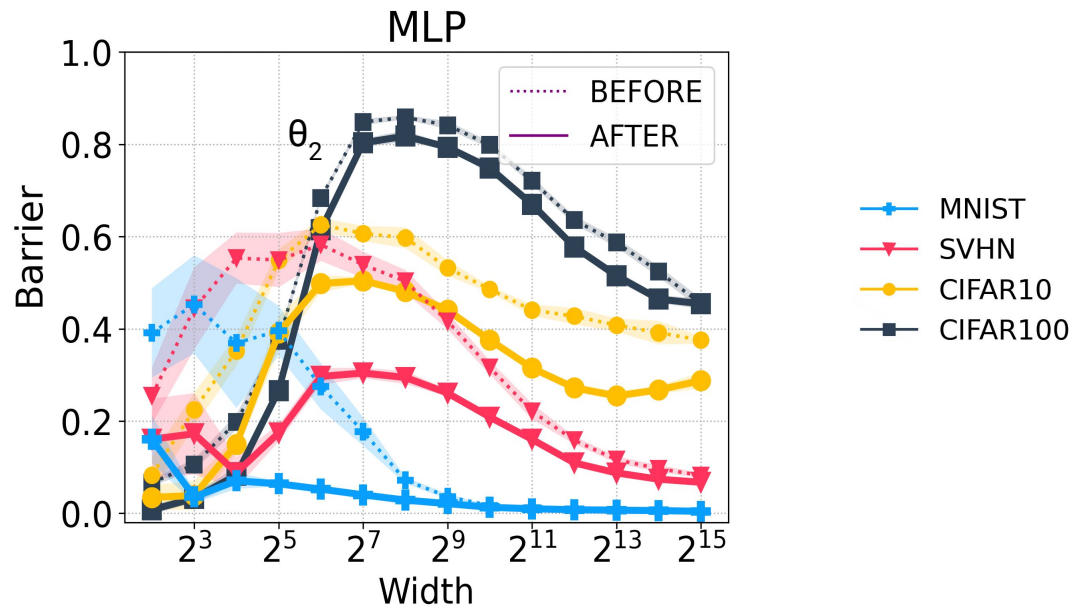
Permutation by Simulated Annealing



Permutation by Simulated Annealing



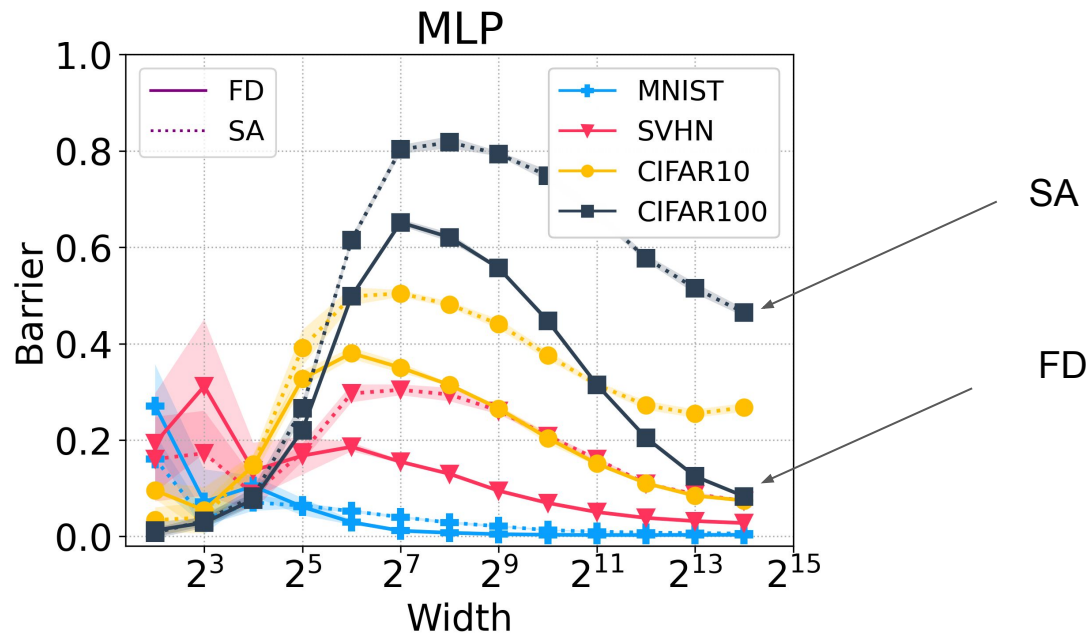
Permutation by Simulated Annealing



Neuron Alignment: **F**unctional **D**ifference

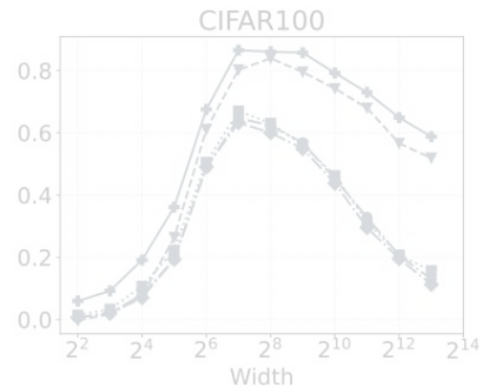
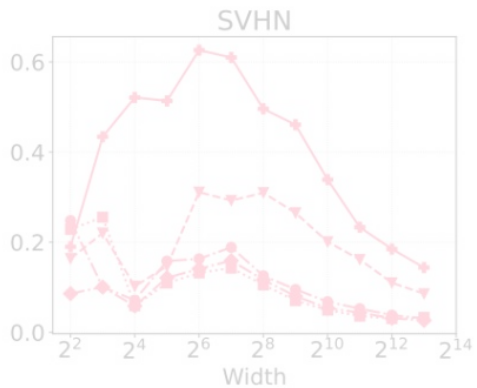
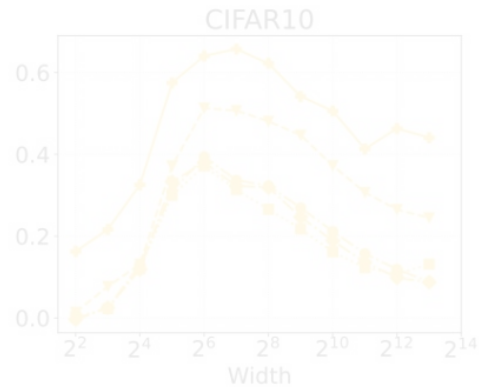
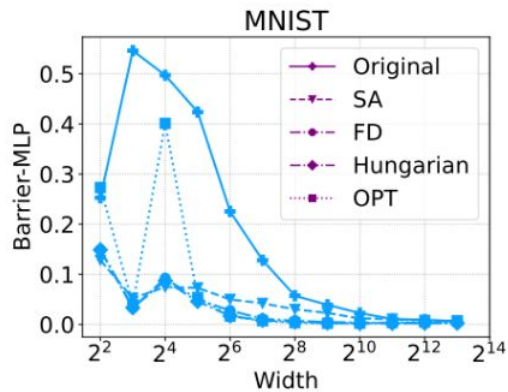
$$\delta E_l^{opt} = \frac{1}{2} (\tilde{\mathbf{w}}_{l,i}^A - \tilde{\mathbf{w}}_{l,j}^B)^\top \cdot \left((\tilde{\mathbf{H}}_{l,i}^A)^{-1} + (\tilde{\mathbf{H}}_{l,j}^B)^{-1} \right)^{-1} \cdot (\tilde{\mathbf{w}}_{l,i}^A - \tilde{\mathbf{w}}_{l,j}^B)$$

Neuron Alignment: Functional Difference

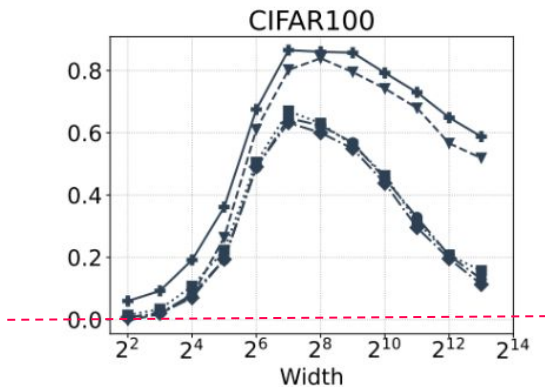
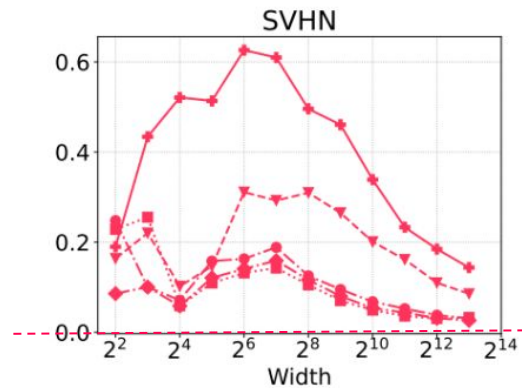
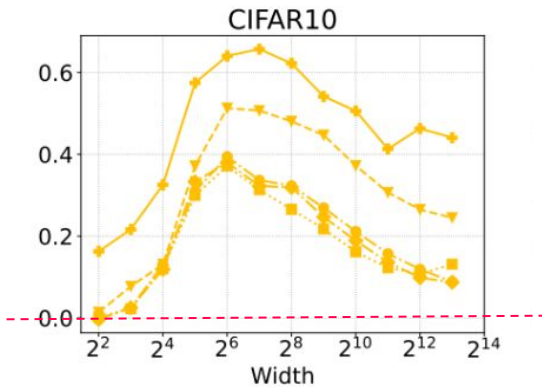
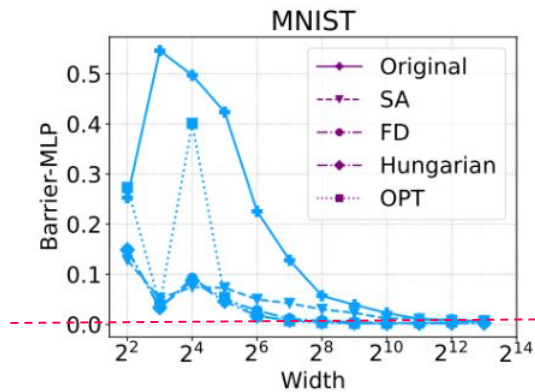


$$\delta E_l^{opt} = \frac{1}{2} (\tilde{\mathbf{w}}_{l,i}^A - \tilde{\mathbf{w}}_{l,j}^B)^\top \cdot \left((\tilde{\mathbf{H}}_{l,i}^A)^{-1} + (\tilde{\mathbf{H}}_{l,j}^B)^{-1} \right)^{-1} \cdot (\tilde{\mathbf{w}}_{l,i}^A - \tilde{\mathbf{w}}_{l,j}^B)$$

Neuron Alignment methods: a comparison

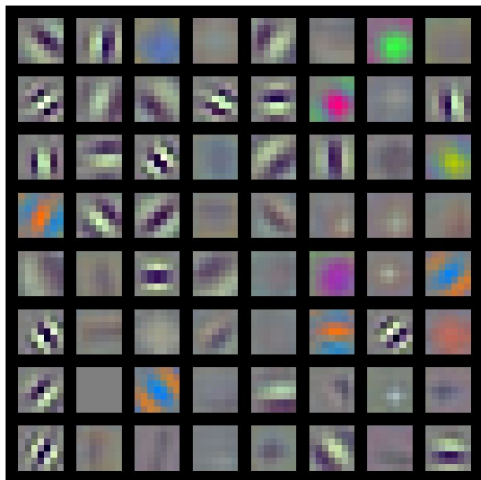


Neuron Alignment methods: a comparison

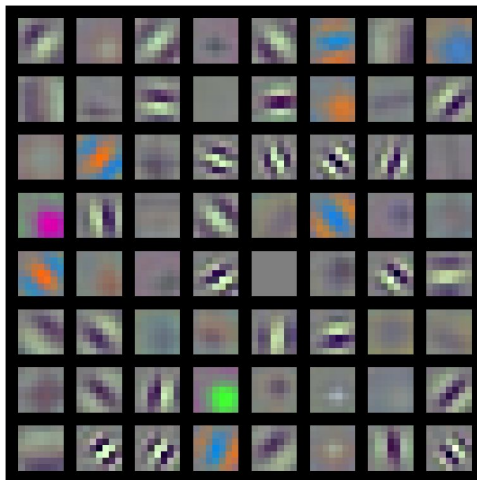


Neuron Alignment: **Correlation Matching**

Network A



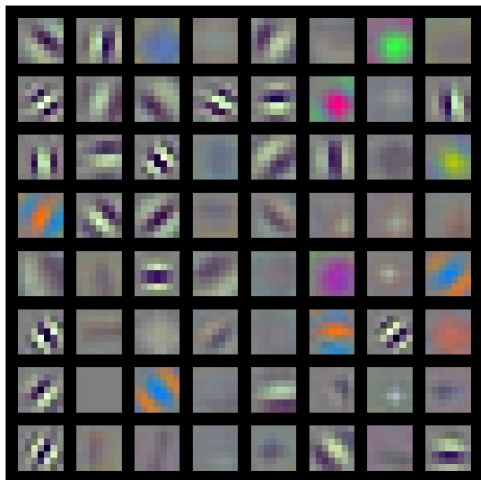
Network B



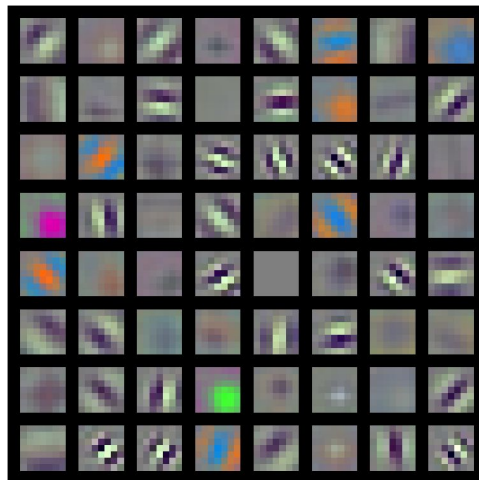
- Resnet-50
- ImageNet
- First layer: 64 filters

Neuron Alignment: Correlation Matching

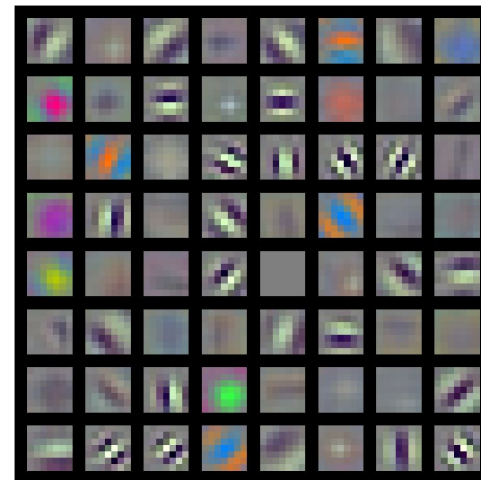
Network A



Network B



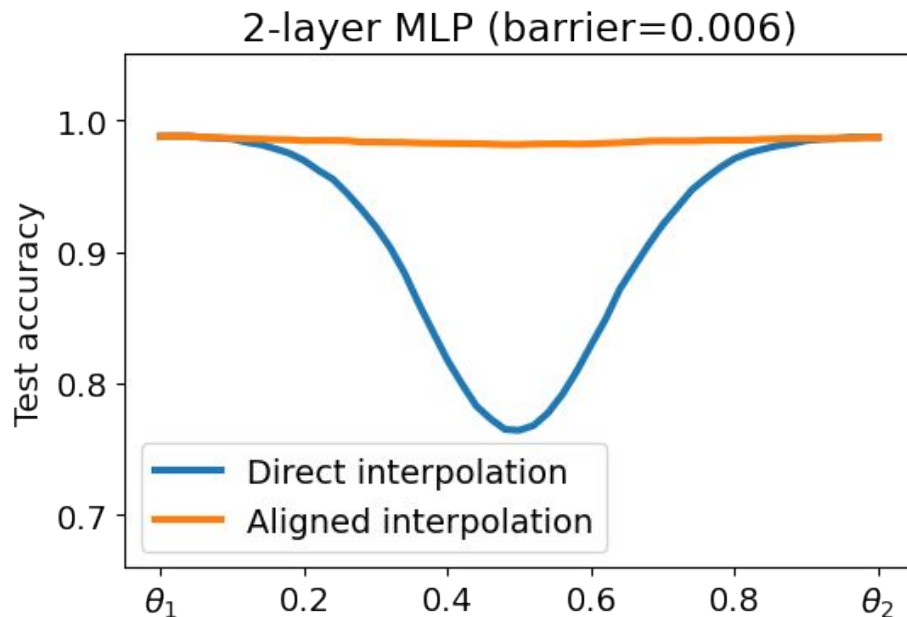
A aligned to B



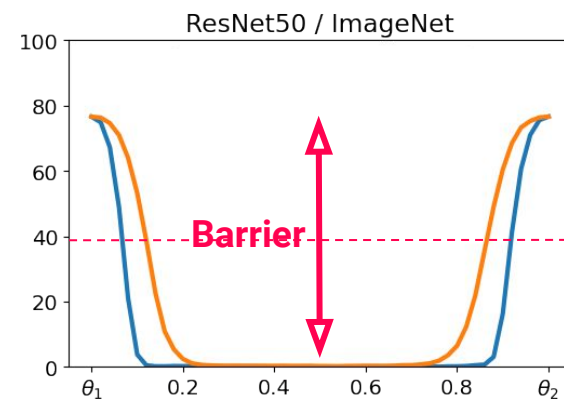
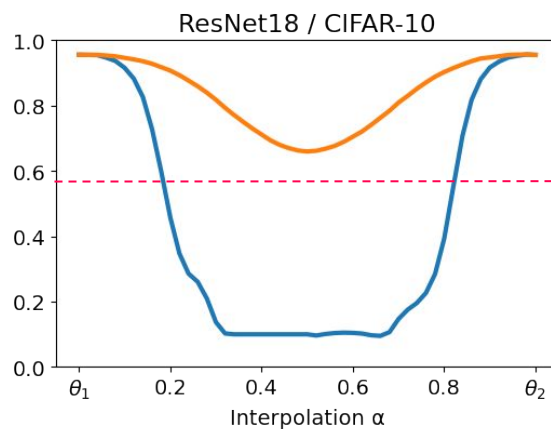
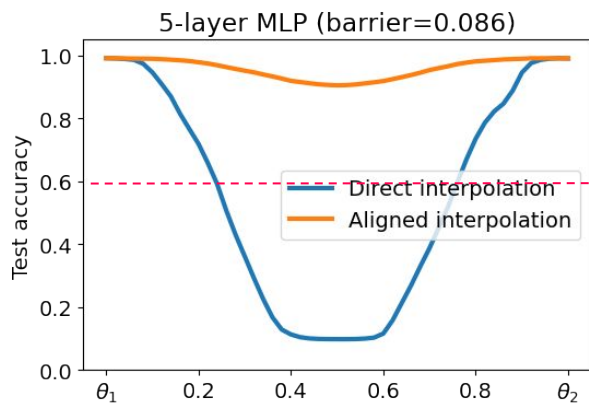
$$\sum_i \text{corr}(X_{l,i}^{(1)}, X_{l,P_l(i)}^{(2)})$$

Neuron Alignment: Correlation Matching

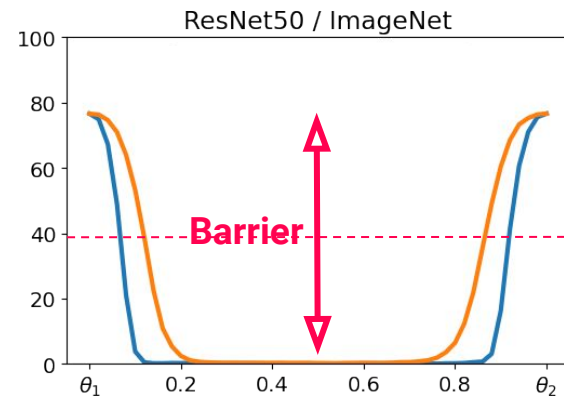
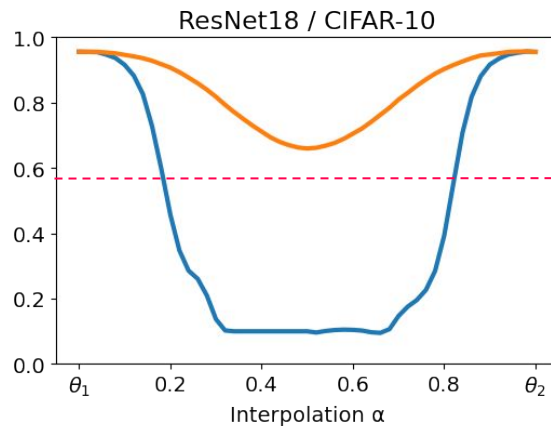
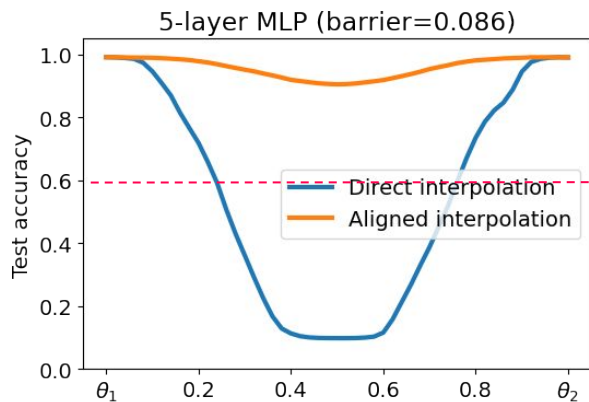
Works for shallow+wide MLPs



Correlation Matching breaks for deeper networks



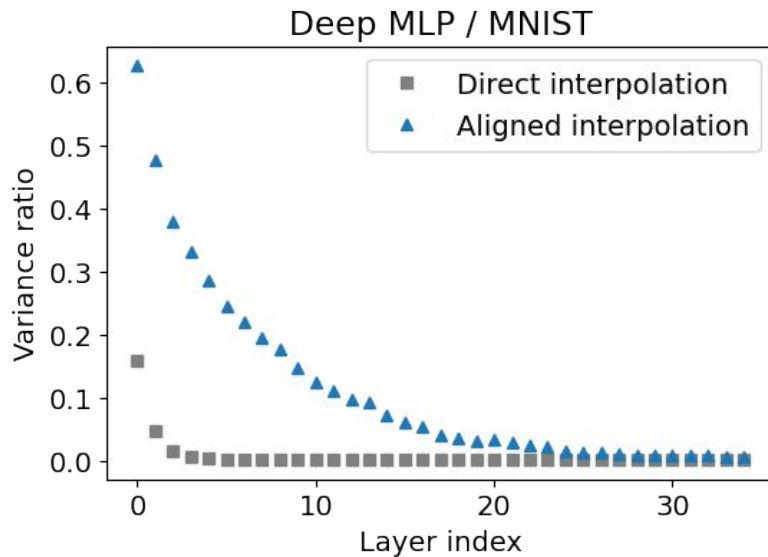
Correlation Matching breaks for deeper networks



But why?

Variance collapse

$$\frac{\sigma_{\alpha}^2}{\frac{\sigma_0^2 + \sigma_1^2}{2}}$$



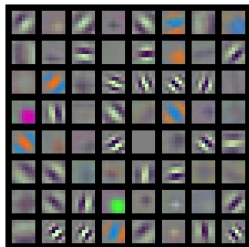
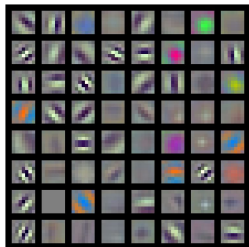
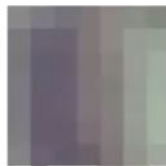
Variance collapse

Filter 9

x_1



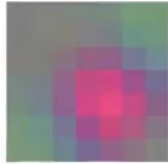
x_2



Variance collapse

Filter 9

x_1



x_2

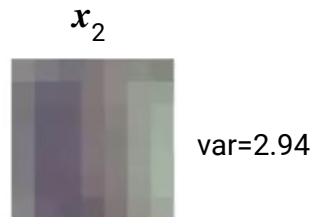
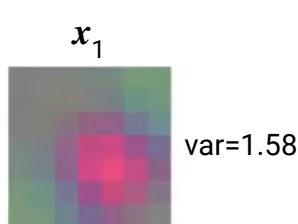


x_α



Variance collapse

Filter 9



$$\begin{aligned}
 \text{Var}(X_\alpha) &= \text{Var}\left(\frac{X_1 + X_2}{2}\right) \\
 &= \frac{\text{Var}(X_1) + \text{Var}(X_2) + 2\text{Cov}(X_1, X_2)}{4} \\
 &= \frac{\text{std}^2(X_1) + \text{std}^2(X_2) + 2 \cdot \text{corr}(X_1, X_2) \cdot \text{std}(X_1)\text{std}(X_2)}{4} \\
 &= \left(\frac{\text{std}(X_1) + \text{std}(X_2)}{2}\right)^2 - \frac{(1 - \text{corr}(X_1, X_2))}{2} \text{std}(X_1)\text{std}(X_2)
 \end{aligned}$$

Variance collapse

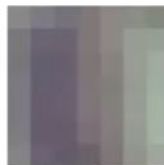
Filter 9

x_1



var=1.58

x_2



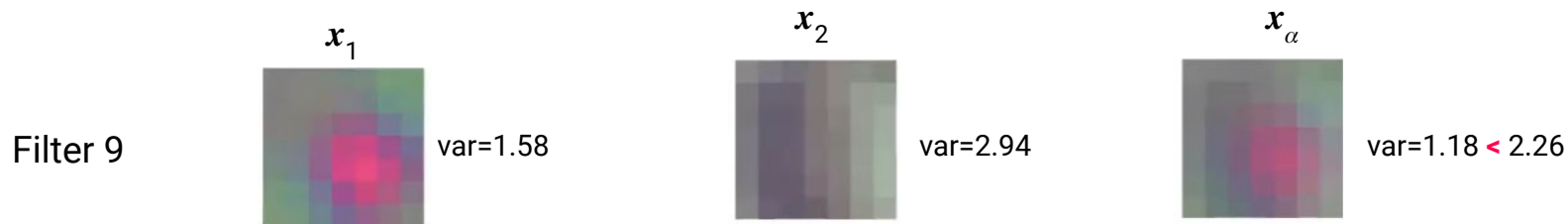
var=2.94

x_α



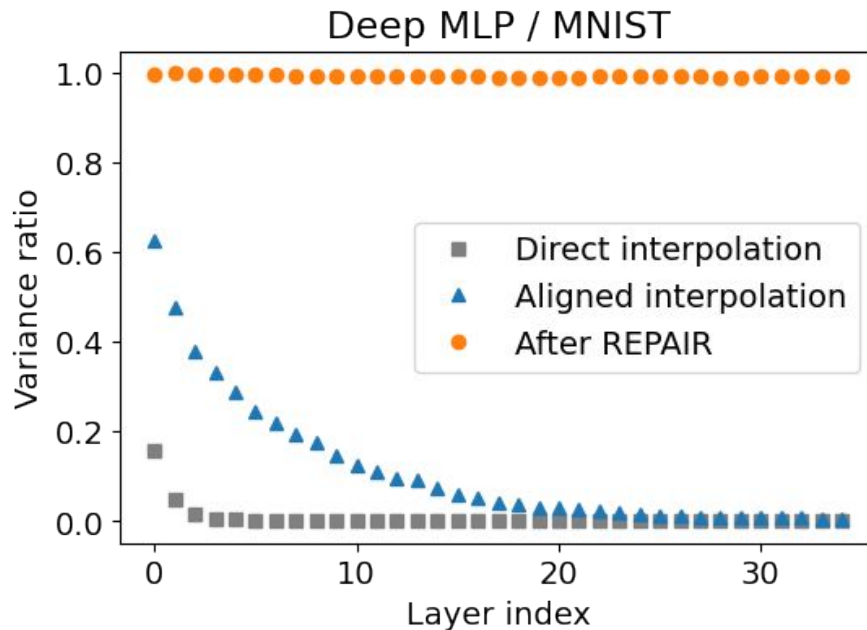
var=1.18 < 2.26

Variance collapse

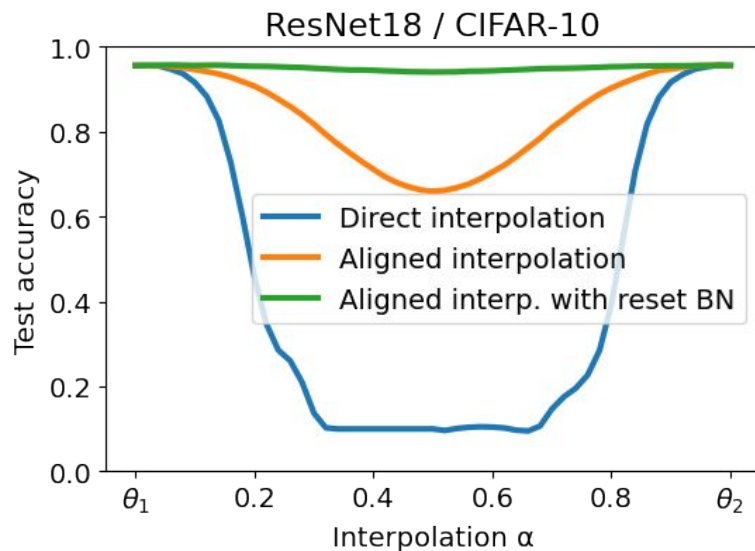


$$\text{Var}(x_\alpha) = \text{Var}\left(\frac{x_1 + x_2}{2}\right) = \frac{1}{4}(\text{std}(x_1)^2 + \text{std}(x_2)^2 + \text{Corr}(x_1, x_2) \cdot \text{std}(x_1) \cdot \text{std}(x_2))$$

REPAIR: Re-estimate Batchnorm statistics



REPAIR: Re-estimate Batchnorm statistics



Part 3: Pre-training Data

Research questions

→ **R1: role of pre-training data**

- ◆ Given a target task, which dataset to pre-train?

→ **R2: role of pre-training method**

- ◆ Given a target task, which pre-train method to choose?
- ◆ supervised ImageNet or contrastive LAION?

Experimental setup

1

Pre-training

CLIP

LAION, YFCC, WIT, Conceptual captions, Redcaps, Shutterstock

2

Finetuning

Few-shot: 1/5/10/20/all samples per class

CIFAR100, DTD, CALTECH101, PETS, REAL (domain net), CLIPART (domain net),
CameraTraps, Cassava Leaf Disease, EuroSAT

Pre-training datasets

LAION



Yellow sandals for women pointy and low heeled Beatnik Françoise Mustard

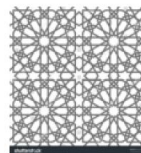


3 Bedrooms Terraced House for sale in Eastbourne Road, Walton, Liverpool, Merseyside, L9



Minimum Wage Barbie

Conceptual captions



Islamic vector geometric ornaments based on traditional arabic art. Oriental seamless pattern. Muslim mosaic. Turkish, Arabian tile on a white background. Mosque ...



Illustration of hand holding the id card. Vector illustration flat design.



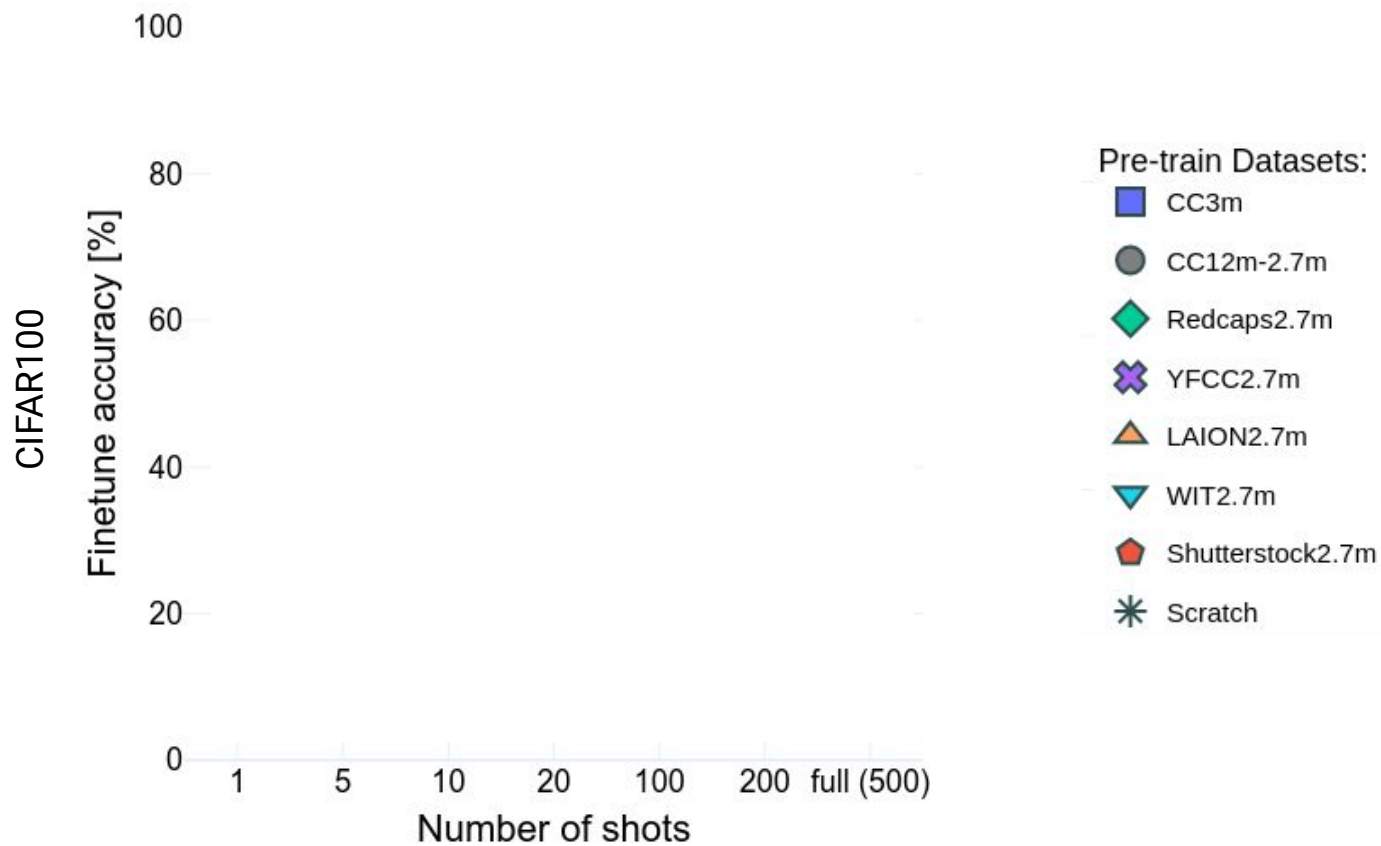
<PERSON>: U. <PERSON> in United States Army. First <PERSON> appointed to that position. First, &, so far, only <PERSON> to serve on Joint Chiefs of Staff. Black H...

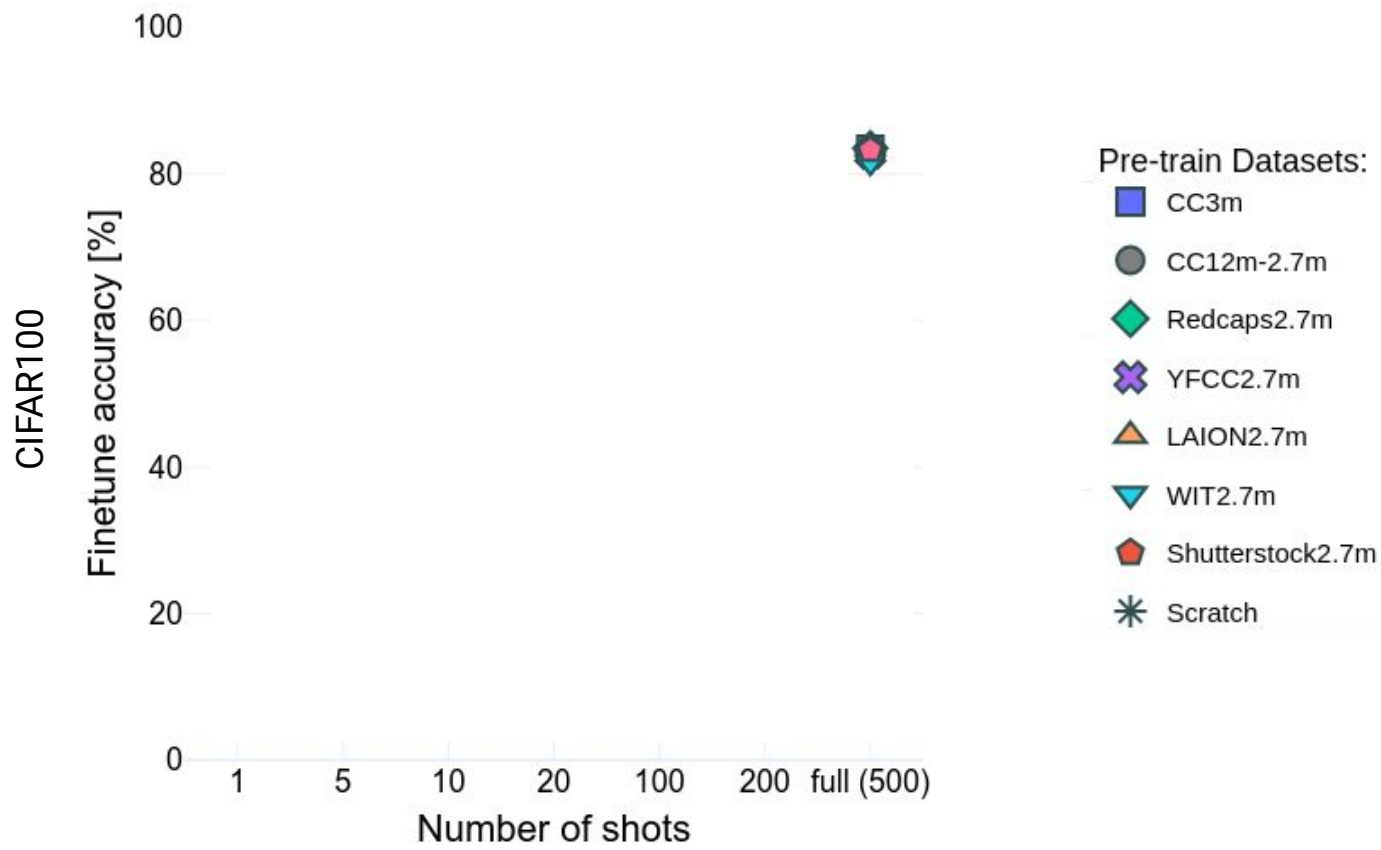
Finetuning datasets

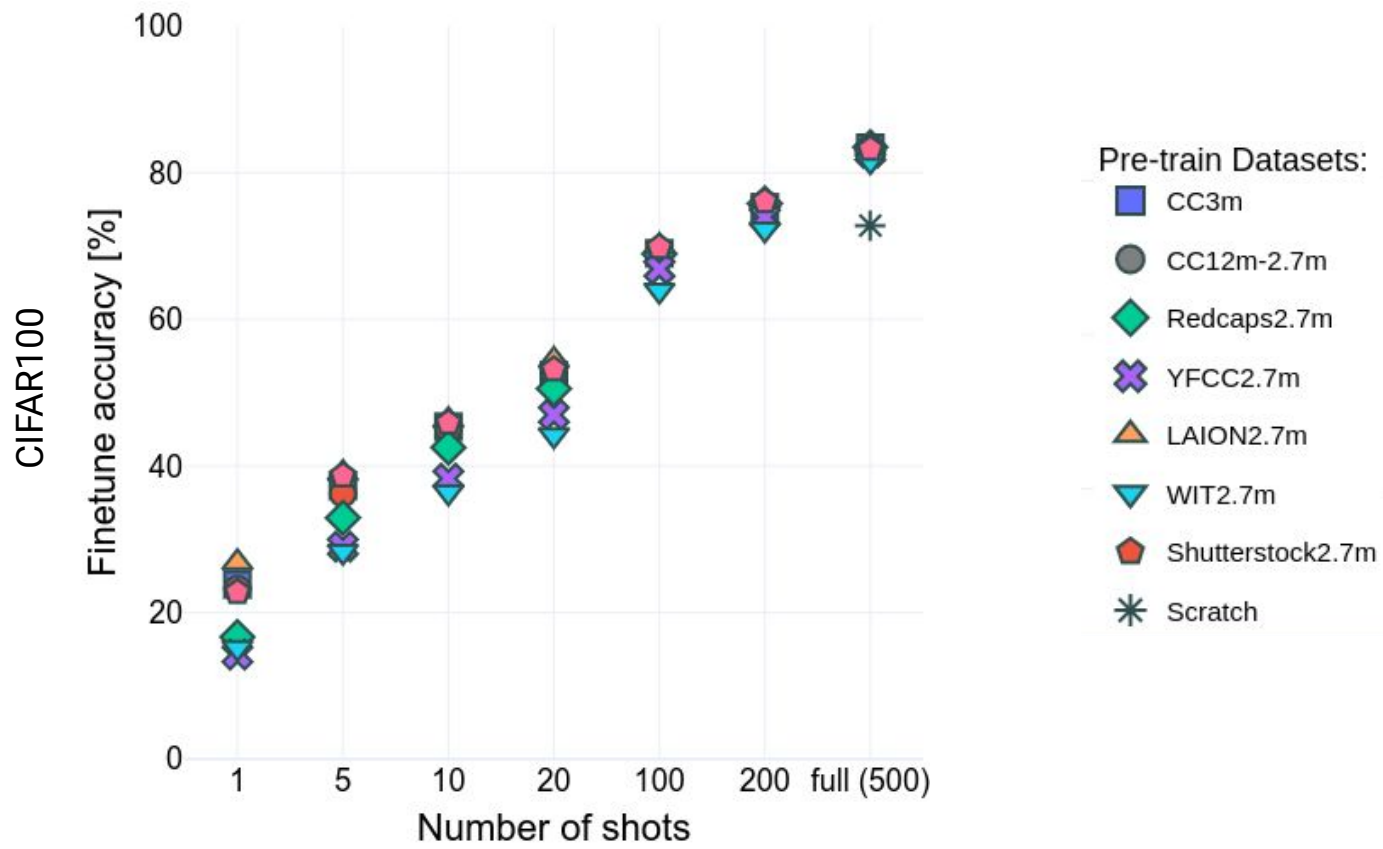
name	CIFAR100	DTD	REAL (domain net)	CLIPART (domain net)	Camera traps	Cassava leaf disease	EuroSAT
samples	50K	5.6K	172K	172K	58K	21K	27K
classes	100	47	345	345	15	5	10



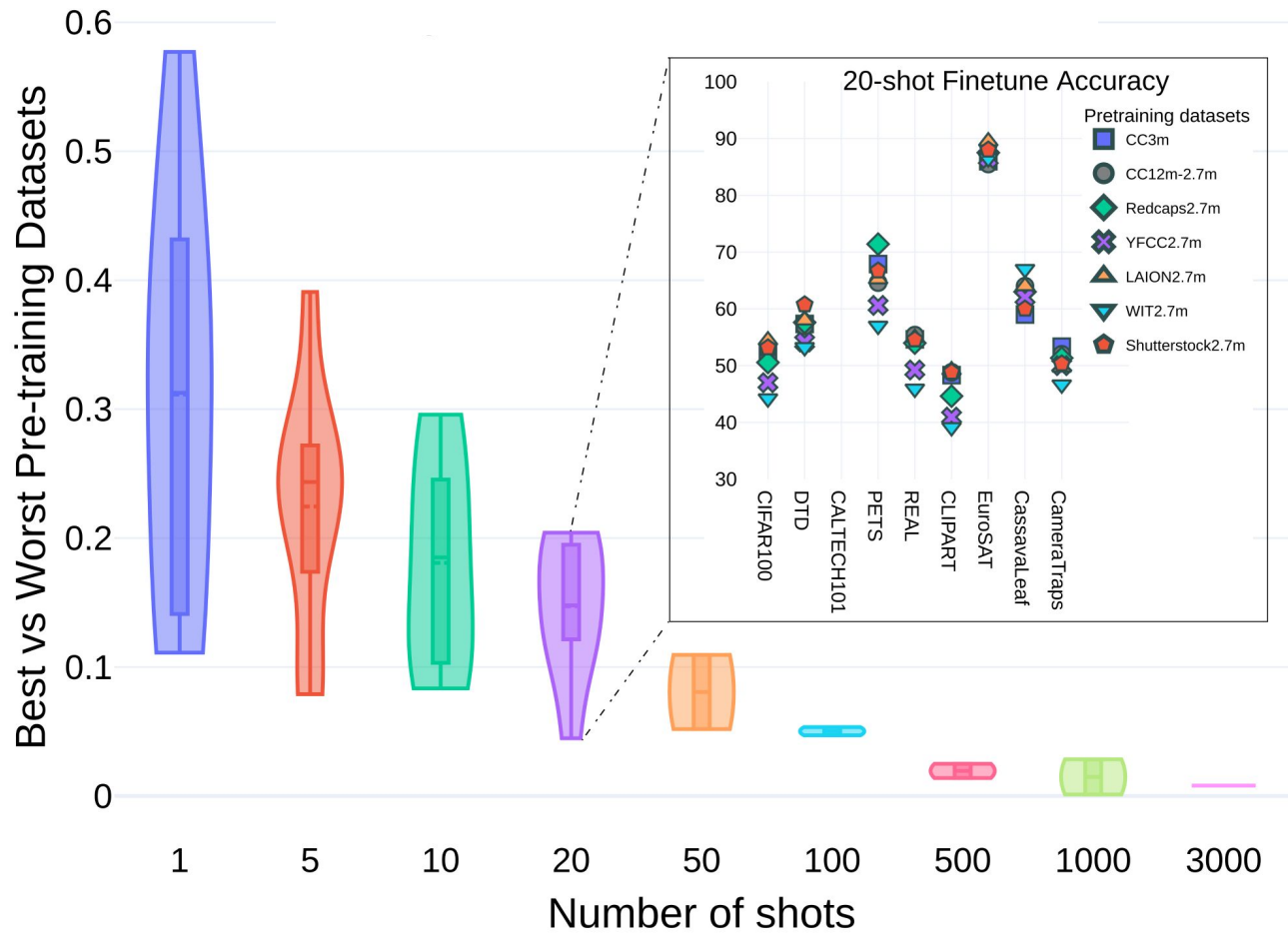
Which dataset to pre-train?



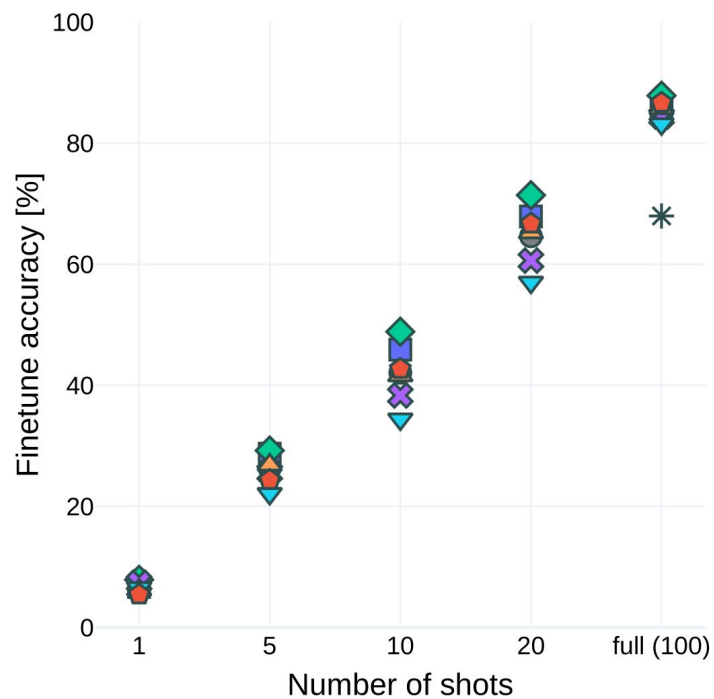




Average over 9 downstream datasets



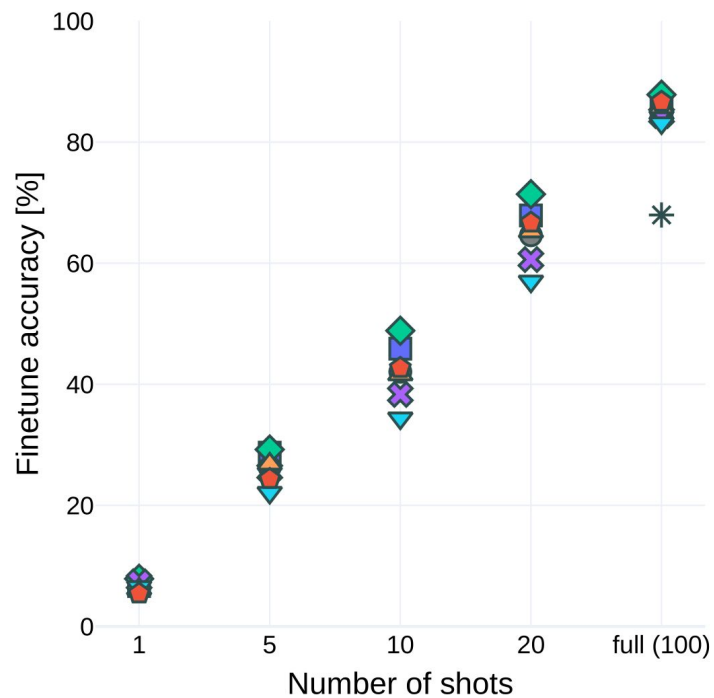
Redcaps on PETS



Pre-train Datasets:

- CC3m
- CC12m-2.7m
- Redcaps2.7m
- YFCC2.7m
- LAION2.7m
- WIT2.7m
- Shutterstock2.7m
- Scratch

Redcaps on PETS



lofoten archipelago by <usr>



bubba is so unbelievably cute when she's sleeping!



the kids got t-shirts



paused the x-men at just the right time.



in a field of yellow and green



eerie section of trail on a long-forgotten country backroad. - long path, catskills park ny



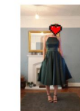
dressing up for the family photo



completed a small remodel of the half bath. first timer.



foggy night in the vancouver forest



duchesse satin wedding guest dress- featuring bonus pockets!



your present condition!



homemade flammkuchen for dinner...



i'm drunk, and this is lucy.



my handsome new neighbour

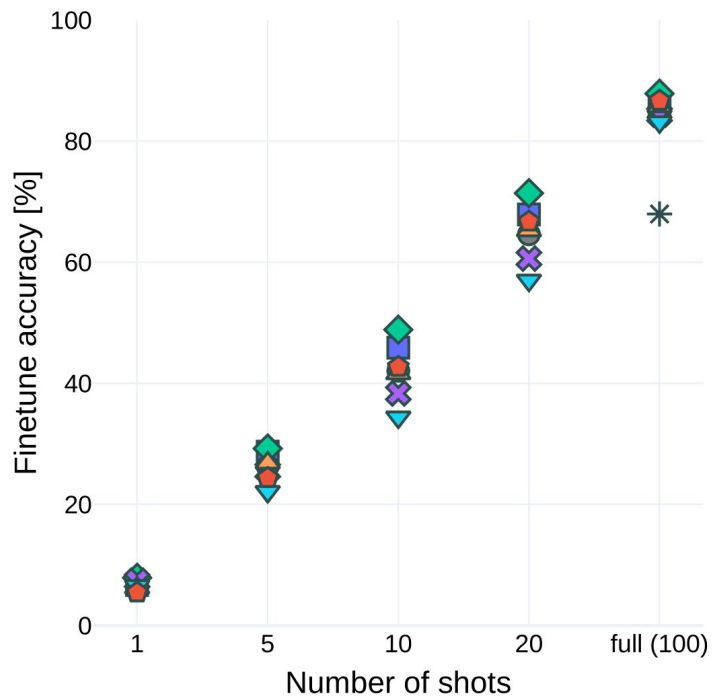


shot from our airbnb porch view on oia on santorini in greece



such a pretty girl

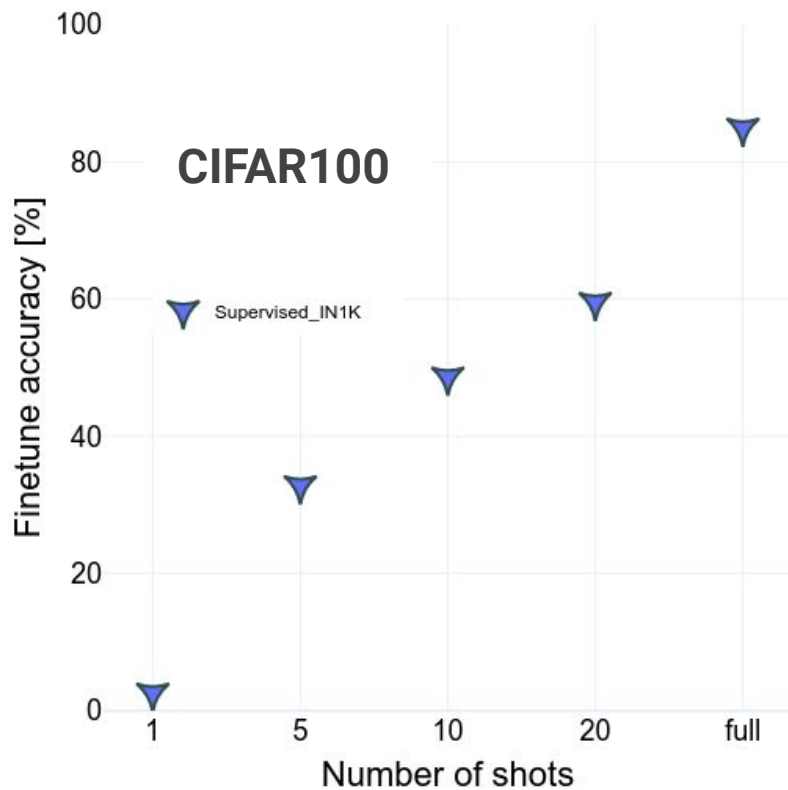
Redcaps on PETS



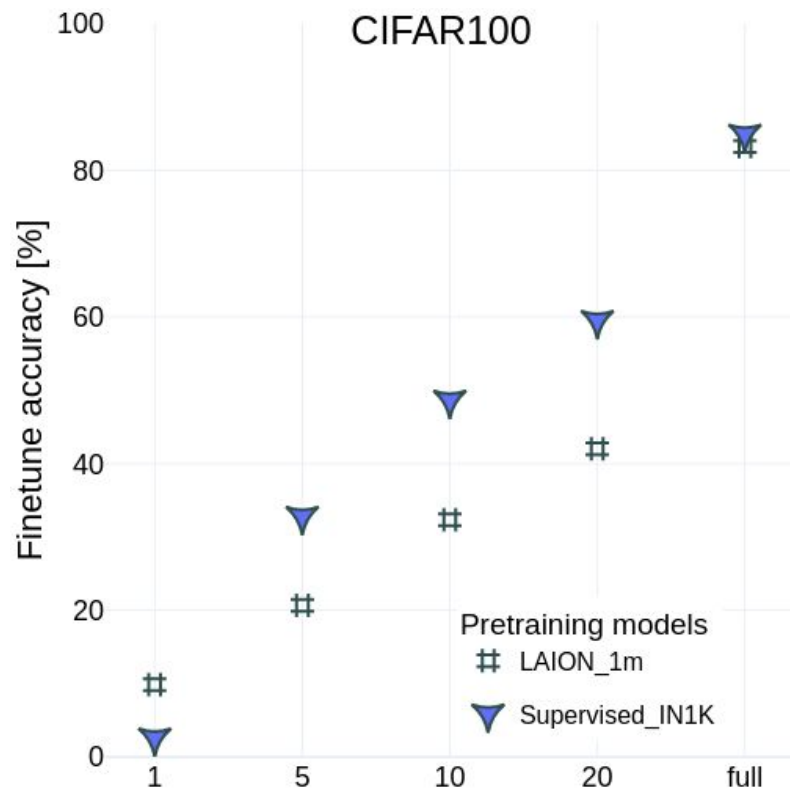
Pre-training dataset	Top 20 words in 1M sample of captions
Shutterstock	background, vector, illustration, design, icon, pattern, texture, style, woman, concept, hand, color, flower, view, template, line, business, logo, card, symbol
Redcaps	day, today, year, time, cat , plant, friend, anyone, picture, baby, guy, week, dog , home, morning, night, month, way, boy, work
YFCC-15m	photo, day, park, street, city, picture, view, time, world, year, house, state, center, part, garden, shot, image, building, road, museum
LAION-15m	photo, stock, image, black, woman, design, set, vector, white, print, home, men, blue, dress, art, card, sale, gold, bag, cover
CC-12m	illustration, stock, art, design, photo, image, background, room, vector, house, home, woman, wedding, style, photography, royalty, car, fashion, girl, world
CC-3m	background, actor, artist, player, illustration, view, woman, man, football, team, tree, premiere, city, vector, day, girl, beach, game, hand, people
WIT	view, church, station, map, house, building, hall, museum, city, location, street, park, river, state, john, county, town, center, bridge, world

Table 2: Most common words in captions of pre-training distributions

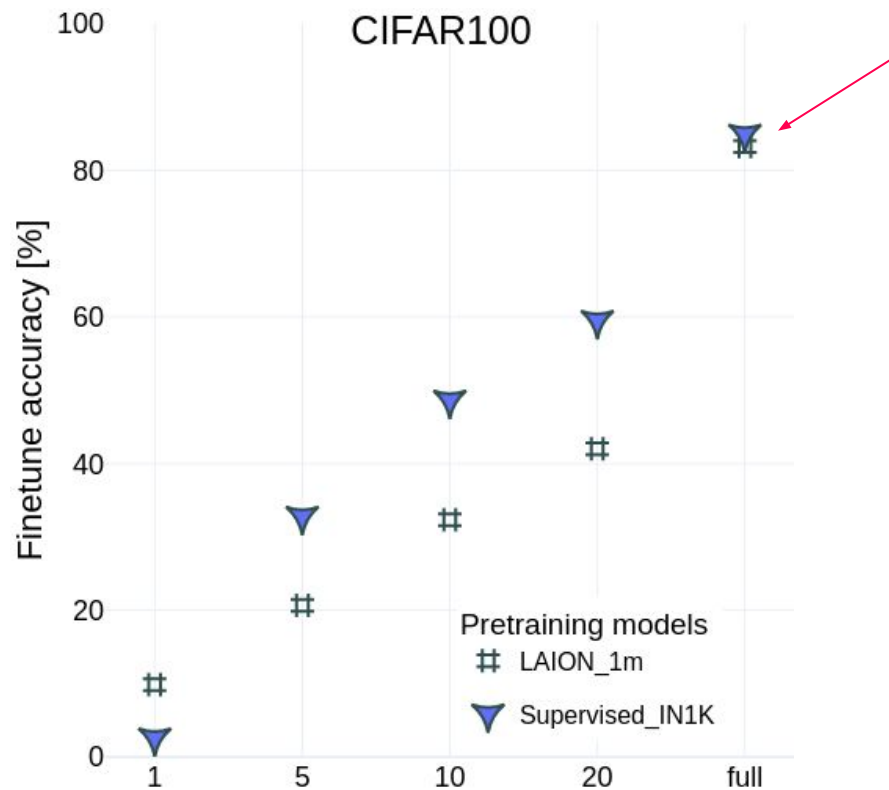
Which pre-train method?



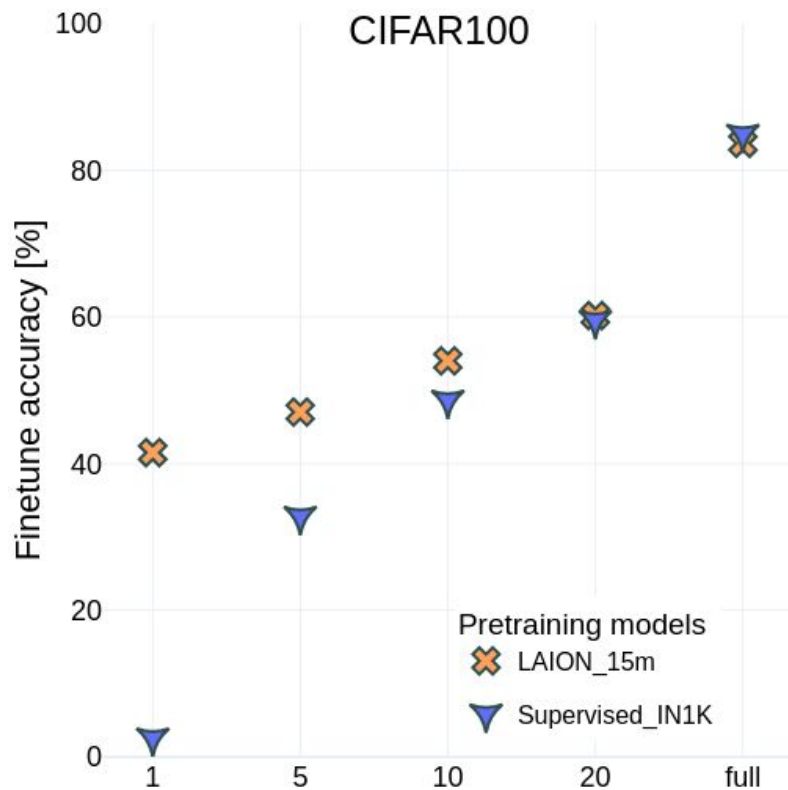
Supervised vs. CLIP



Supervised vs. CLIP



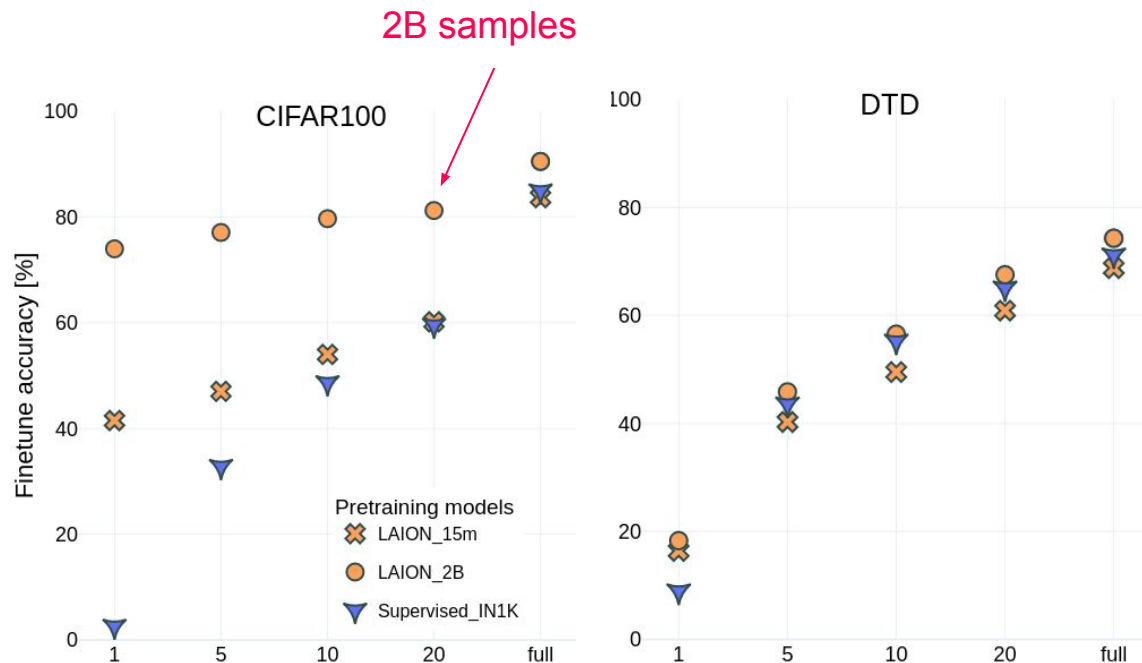
Adding 15x more data?



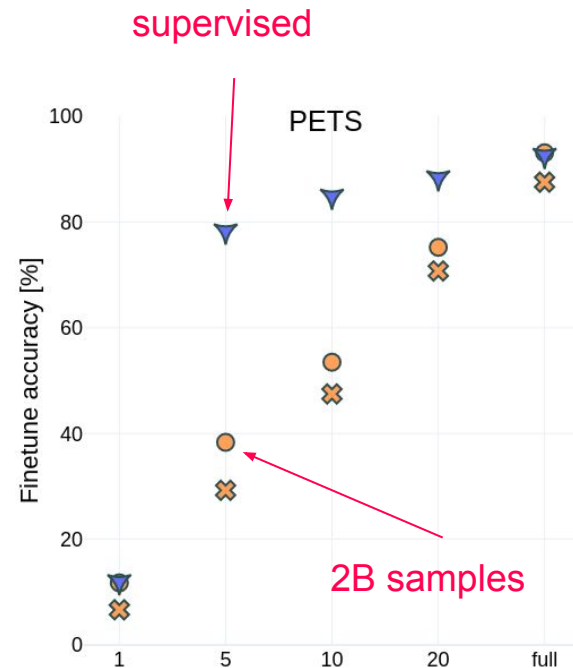
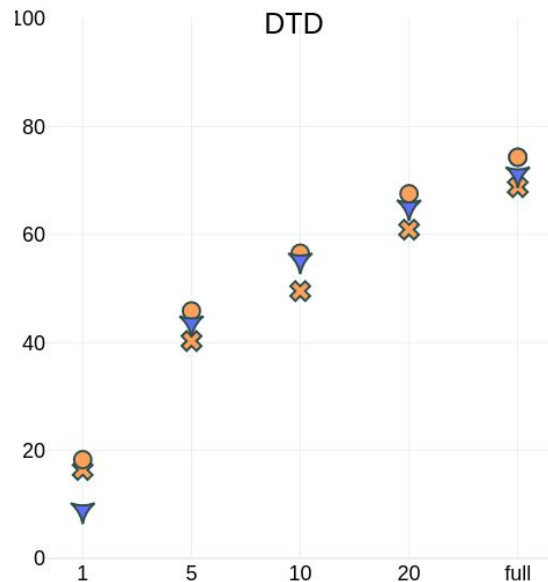
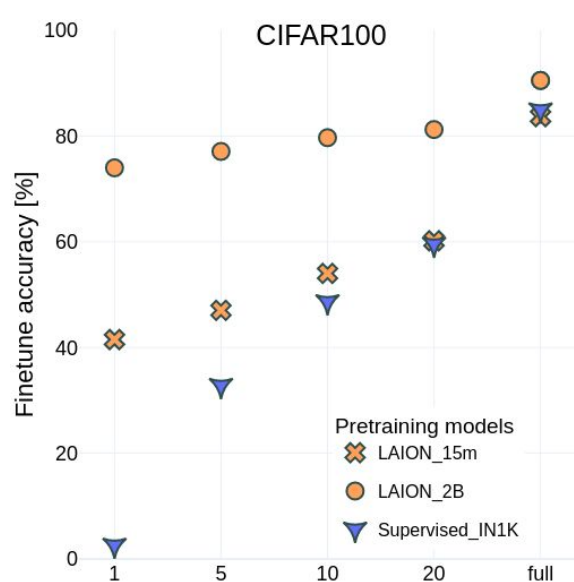
Adding 15x more data?



What if we scale to 2B samples?



What if we scale to 2B samples?



Take away

- Sparsity:
 - There are several motivations for sparsity, one is improving generalization.
 - Sparsity has different effects when combined with supervised and semi-supervised training.

Take away

- Loss landscape:
 - Studying the loss landscape of neural networks has implications on model generalization.
 - Accounting for permutation invariance, barriers can be eliminated.
 - New lens to loss landscape: we took the first steps towards understanding ensembles and distributed training.

Take away

- Role of data:
 - Changing the pre-training dataset leads to noticeable differences in few-shot transfer performance.
 - Specific datasets like shutterstock perform well on almost all studied target tasks.
 - Data curation matters. We need 15-2000X more data to compensate for labeling.

Thanks for your attention