Robust fine-tuning of zero-shot models













Mitchell Wortsman*, Gabriel Ilharco*, Jong Wook Kim, Mike Li, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, Ludwig Schmidt

























Train





















Test















Test































Evaluating models under natural distribution shift

Measuring robustness to natural distribution shift. Taori et al., 2020

) 3

Evaluating models under natural distribution shift

Distribution *D*

ImageNet (Deng et al.)







Measuring robustness to natural distribution shift. Taori et al., 2020







) 3

Evaluating models under natural distribution shift

Distribution *D*





Measuring robustness to natural distribut





ObjectNet (Barbu et al.)

ImageNetV2 (Recht et al.)

ImageNet-A (Hendrycks et al.)







































ImageNet

Source: https://openai.com/blog/clip/

IMAGENET RESNET101

76.2%

5





ImageNet



ObjectNet

Source: https://openai.com/blog/clip/

IMAGENET RESNET101

76.2%

32.6%



CLIP: Connecting Text and Images

We're introducing a neural network called CLIP which efficiently learns visual concepts from natural language supervision. CLIP can be applied to any visual classification benchmark by simply providing the names of the visual categories to be recognized, similar to the "zero-shot" capabilities of GPT-2 and GPT-3.

January 5, 2021 15 minute read



CLIP is trained to associate **image-caption pairs** on the internet, and can perform zero-shot inference.

CLIP is trained to associate **image-caption pairs** on the internet, and can perform zero-shot inference.



✓ a photo of **guacamole**, a type of food. × a photo of **ceviche**, a type of food. × a photo of **edamame**, a type of food. × a photo of **tuna tartare**, a type of food. × a photo of **hummus**, a type of food.



ImageNet

IMAGENET RESNET101

CLIP VIT-L

76.2%

76.2%





ImageNet



ObjectNet

IMAGENET **RESNET101** CLIP VIT-L 76.2% 76.2% 32.6% 72.3%





ImageNet



ImageNet V2



ImageNet Rendition



ObjectNet



ImageNet Sketch



ImageNet Adversarial

IMAGENET RESNET101	CLIP VIT-L
76.2%	76.2%
64.3%	70.1%
37.7%	88.9%
32.6%	72.3%
25.2%	60.2%
2.7%	77.1%



Source: <u>https://paperswithcode.com/sota/image-classification-on-imagenet</u>

To improve a models accuracy it is *fine-tuned* on data from the domain of interest.

interest.

But: despite accuracy improvements, *fine-tuning* CLIP deteriorates robustness.

To improve a models accuracy it is *fine-tuned* on data from the domain of

interest.

But: despite accuracy improvements, *fine-tuning* CLIP deteriorates robustness.

This work:

To improve a models accuracy it is *fine-tuned* on data from the domain of

interest.

But: despite accuracy improvements, *fine-tuning* CLIP deteriorates robustness.

This work:

To improve a models accuracy it is *fine-tuned* on data from the domain of

Can we fine-tune CLIP while maintaining robustness?



Measuring robustness to natural distribution shift. Taori et al., 2020 12

Accuracy on ${\cal D}$



Measuring robustness to natural distribution shift. Taori et al., 2020

13



Evaluating Machine Accuracy on ImageNet. Shankar et al., 2020



Measuring robustness to natural distribution shift. Taori et al., 2020



Measuring robustness to natural distribution shift. Taori et al., 2020



CLIP. Radford et al., 2020

Accuracy on ${\cal D}$



CLIP. Radford et al., 2020

Accuracy on ${\cal D}$



CLIP. Radford et al., 2020

Accuracy on ${\cal D}$


Accuracy on ${\cal D}$



Accuracy on \mathcal{D}



Accuracy on ${\cal D}$













Weight-space ensembling = linearly interpolating in weight space



Weight-space ensembling = linearly interpolating in weight space









CLIP zero-shot

- Linear fit (CLIP zero-shot)
- CLIP fine-tuned end-to-end
- Weight-space ensemble (end-to-end)
- Best OOD without reducing ID
- Standard ImageNet models
 - Linear fit (standard ImageNet models)

$$y = x$$







- Baselines
- Choosing the mixing coefficient
- More distribution shifts and datasets
- Why does this work

- Baselines
- Choosing the mixing coefficient
- More distribution shifts and datasets
- Why does this work













- Baselines
- Choosing the mixing coefficient
- More distribution shifts and datasets
- Why does this work



ImageNet (Deng et al.)



ImageNetV2 (Recht et al.)



ImageNet-R (Hendrycks et al.)















ObjectNet (Barbu et al.)



ImageNet-A (Hendrycks et al.)























ImageNet (Deng et al.)



ImageNetV2 (Recht et al.)



ImageNet-R (Hendrycks et al.)















ObjectNet (Barbu et al.)



ImageNet-A (Hendrycks et al.)























- Baselines
- Choosing the mixing coefficient
- More distribution shifts and datasets
- Why does this work





iWildCam



Beery et al., 2018



iWildCam



Beery et al., 2018



iWildCam



Beery et al., 2018

FMoW

Christie et al., 2018

	Train			Test	
Satellite Image (x)					
Year / Region (d)	2002 / Americas	2009 / Africa	2012 / Europe	2016 / Americas	2017 / Africa
Building / Land Type (y)	shopping mall	multi-unit residential	road bridge	recreational facility	educational institution







iWildCam



Beery et al., 2018

FMoW

+3.7pp OOD

Christie et al., 2018

	Train			Test	
Satellite Image (x)					
Year / Region (d)	2002 / Americas	2009 / Africa	2012 / Europe	2016 / Americas	2017 / Africa
Building / Land Type (y)	shopping mall	multi-unit residential	road bridge	recreational facility	educational institution







iWildCam



Beery et al., 2018

FMoW

+3.7pp OOD

Christie et al., 2018

	Train			Test	
Satellite Image (x)					
Year / Region (d)	2002 / Americas	2009 / Africa	2012 / Europe	2016 / Americas	2017 / Africa
Building / Land Type (y)	shopping mall	multi-unit residential	road bridge	recreational facility	educational institution





CIFAR-10.1. Recht et al., 2019

CIFAR-10.2. Lu et al., 2020







iWildCam



Beery et al., 2018

FMoW

+3.7pp OOD

Christie et al., 2018

	Train			Test	
Satellite Image (x)					
Year / Region (d)	2002 / Americas	2009 / Africa	2012 / Europe	2016 / Americas	2017 / Africa
Building / Land Type (y)	shopping mall	multi-unit residential	road bridge	recreational facility	educational institution

+2.2pp OOD



+3.0pp OOD



CIFAR-10.1. Recht et al., 2019

CIFAR-10.2. Lu et al., 2020







iWildCam



Beery et al., 2018

FMoW

+3.7pp OOD

Christie et al., 2018

	Train			Test	
Satellite Image (x)					
Year / Region (d)	2002 / Americas	2009 / Africa	2012 / Europe	2016 / Americas	2017 / Africa
Building / Land Type (y)	shopping mall	multi-unit residential	road bridge	recreational facility	educational institution

+2.2pp OOD



Predicted: domestic cat



+3.0pp OOD



CIFAR-10.1. Recht et al., 2019

CIFAR-10.2. Lu et al., 2020

Predicted: monkey



ImageNet-Vid-Robust Shankar et al., 2019 **YTBBRobust**









iWildCam



Beery et al., 2018

FMoW

+3.7pp OOD

Christie et al., 2018

	Train			Test	
Satellite Image (x)					
Year / Region (d)	2002 / Americas	2009 / Africa	2012 / Europe	2016 / Americas	2017 / Africa
Building / Land Type (y)	shopping mall	multi-unit residential	road bridge	recreational facility	educational institution

+2.2pp OOD



Predicted: domestic cat



+3.0pp OOD



Predicted: monkey



CIFAR-10.1. Recht et al., 2019

CIFAR-10.2. Lu et al., 2020

+8.3pp OOD

ImageNet-Vid-Robust

Shankar et al., 2019

YTBBRobust +14.7pp OOD













ImageNet (Deng et al., 2009)





Describable Textures (Cimpoi et al., 2014)









/s/subway_station/platform (337)







/c/control_tower/outdoor (110)



/b/bazaar/indoor (44)



/d/driveway (138)



/a/auto_factory (23)


ImageNet (Deng et al., 2009)





Describable Textures (Cimpoi et al., 2014)



+3.3pp ID

CIFAR-10 & CIFAR-100 (Krizhevsky et al., 2009)







/s/subway_station/platform (337)









+2.0pp ID

(Bossard et al., 2014)



+1.8pp ID

/c/control_tower/outdoor (110)

<u>et</u> al., 2016)



/b/bazaar/indoor (44)

/a/auto_factory (23)



/d/driveway (138)





- Baselines
- Choosing the mixing coefficient
- More distribution shifts and datasets
- Why does this work

Outline







Observation 1: (Frankle et al., 2020)



Observation 1: (Frankle et al., 2020)

Ensembling in weight space — making predictions via

$$f\left(x,\frac{1}{2}\left(\theta_1^T+\theta_2^T\right)\right)$$

fails, performing no better than random chance.



Observation 2: (Frankle et al., 2020)



Observation 2 (Frankle et al., 2020)

When part of the training trajectory is shared, there exists a *linear* path in weight space between the two solutions along which loss remains low.



Observation 3 (Izmailov et al., 2018)

Weights found through SGD tend to lie at the periphery of a large, flat minimum.





Observation 3 (Izmailov et al., 2018)

Weights found through SGD tend to lie at the periphery of a large, flat minimum.

To get closer to the center, Stochastic Weight Averaging (SWA) bounces around the minimum while saving checkpoints and returns the average.





Observation 4 (Neyshabur, Sedghi, Zhang, 2020)

Linear mode connectivity between fine-tuned solutions.



Observation 5

Linear mode connectivity between the zero-shot and fine-tuned solution on both the OOD and ID dataset.



$\mathsf{Acc}((1 - \alpha) \cdot \theta_0 + \alpha \cdot \theta_1) \ge (1 - \alpha) \cdot \mathsf{Acc}(\theta_0) + \alpha \cdot \mathsf{Acc}(\theta_1)$

 $\operatorname{Acc}((1 - \alpha) \cdot \theta_0 + \alpha \cdot \theta_1) \ge (1 - \alpha) \cdot \operatorname{Acc}(\theta_0) + \alpha \cdot \operatorname{Acc}(\theta_1)$







Schematic: average test error on all datasets

CLIP zero-shot (high ID error, low OOD error)

CLIP fine-tuned end-to-end (low ID error, high OOD error)

Weight-space ensemble (low ID and OOD error)

Robustness to distribution shift is an important step in ML.

Robustness to distribution shift is an important step in ML.

- Large-scale pre-trained models such as CLIP are a promising direction.

Robustness to distribution shift is an important step in ML.

robustness.

- Large-scale pre-trained models such as CLIP are a promising direction.
- Standard fine-tuning improves accuracy in-distribution, but deteriorates

Robustness to distribution shift is an important step in ML.

robustness.

WiSE-FT mitigates the compromise between high accuracy and robustness with no extra compute during fine-tuning or inference.

- Large-scale pre-trained models such as CLIP are a promising direction.
- Standard fine-tuning improves accuracy in-distribution, but deteriorates