

# Frozen Feature Augmentation for Few-Shot Image Classification



Andreas Bär<sup>+</sup>  
TU Braunschweig  
Germany



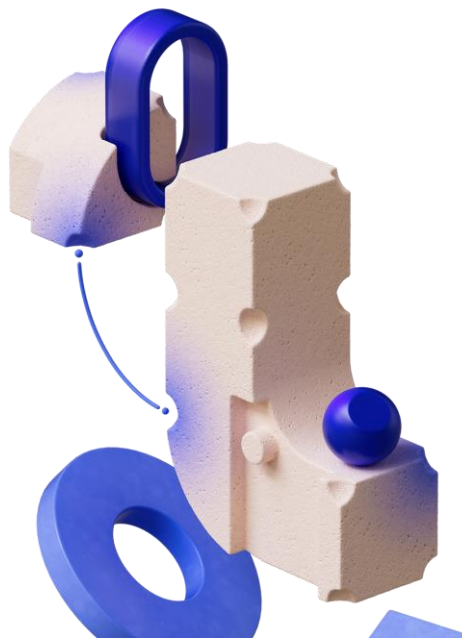
Neil Houlsby  
Google DeepMind  
Switzerland



Mostafa Dehghani  
Google DeepMind  
USA



Manoj Kumar §  
Google DeepMind  
Netherlands



§ Project lead.  
+ Work done during summer internship at Google DeepMind.

# Introduction

## From Pretraining to Transfer Learning With Frozen Features

Vision models are usually pretrained on **large-datasets**, e.g., ImageNet-21k [Deng+, 2009] or JFT [Zhai+, 2022], and then adapted.

[Deng et al., „ImageNet: A Large-Scale Hierarchical Image Database“, CVPR, June 2009]

[Zhai et al., „Scaling Vision Transformers“, CVPR, June 2022]

**Linear probing** [Radford+, 2021] is an **effective** method to **transfer** vision models to other tasks [Dehghani+, 2023] using frozen features.

[Radford et al., „Learning Transferable Visual Models From Natural Language Supervision“, ICML, July 2021]

[Dehghani et al., „Scaling Vision Transformers to 22 Billion Parameters“, ICML, July 2023]

# Introduction

## From Pretraining to Transfer Learning With Frozen Features

Vision models are usually pretrained on **large-datasets**, e.g., ImageNet-21k [Deng+, 2009] or JFT [Zhai+, 2022], and then adapted.

[Deng et al., „ImageNet: A Large-Scale Hierarchical Image Database“, CVPR, June 2009]

[Zhai et al., „Scaling Vision Transformers“, CVPR, June 2022]

**Linear probing** [Radford+, 2021] is an **effective** method to **transfer** vision models to other tasks [Dehghani+, 2023] using frozen features.

[Radford et al., „Learning Transferable Visual Models From Natural Language Supervision“, ICML, July 2021]

[Dehghani et al., „Scaling Vision Transformers to 22 Billion Parameters“, ICML, July 2023]

**Known methods** [Radford+, 2021; Dehghani+, 2023; Zhai+, 2023] **do not employ data augmentation** when training on frozen features.

[Zhai et al., „Sigmoid Loss for Language-Image Pretraining“, ICCV, October 2023]

**At the same time, data augmentation is known to be effective** [Cubuk+, 2019; Hendrycks+, 2020; Cubuk+, 2020; Müller & Hutter, 2021].

[Cubuk et al., „AutoAugment: Learning Augmentation Strategies From Data“, CVPR, June 2019]

[Hendrycks et al., „AugMix: A Simple Data Processing Method to Improve Robustness and Uncertainty“, ICLR, April 2020]

[Cubuk et al., „RandAugment: Practical Automated Data Augmentation With a Reduced Search Space“, NeurIPS, December 2020]

[Müller and Hutter, „TrivialAugment: Tuning-Free Yet State-of-the-Art Data Augmentation“, ICCV, October 2021]

# Introduction

## From Pretraining to Transfer Learning With Frozen Features

Vision models are usually pretrained on **large-datasets**, e.g., ImageNet-21k [Deng+, 2009] or JFT [Zhai+, 2022], and then adapted.

[Deng et al., „ImageNet: A Large-Scale Hierarchical Image Database“, CVPR, June 2009]

[Zhai et al., „Scaling Vision Transformers“, CVPR, June 2022]

**Linear probing** [Radford+, 2021] is an **effective** method to **transfer** vision models to other tasks [Dehghani+, 2023] using frozen features.

[Radford et al., „Learning Transferable Visual Models From Natural Language Supervision“, ICML, July 2021]

[Dehghani et al., „Scaling Vision Transformers to 22 Billion Parameters“, ICML, July 2023]

Known methods [Radford+, 2021; Dehghani+, 2023; Zhai+, 2023] **do not employ data augmentation** when training on frozen features.

[Zhai et al., „Sigmoid Loss for Language-Image Pretraining“, ICCV, October 2023]

At the same time, **data augmentation is known to be effective** [Cubuk+, 2019; Hendrycks+, 2020; Cubuk+, 2020; Müller & Hutter, 2021].

[Cubuk et al., „AutoAugment: Learning Augmentation Strategies From Data“, CVPR, June 2019]

[Hendrycks et al., „AugMix: A Simple Data Processing Method to Improve Robustness and Uncertainty“, ICLR, April 2020]

[Cubuk et al., „RandAugment: Practical Automated Data Augmentation With a Reduced Search Space“, NeurIPS, December 2020]

[Müller and Hutter, „TrivialAugment: Tuning-Free Yet State-of-the-Art Data Augmentation“, ICCV, October 2021]

In a **data-constrained, few-shot setup** we investigate:

### Research question:

Can we effectively combine image augmentations and frozen features to train a lightweight model?

# Method Description

## Training on Top of Frozen Features in Three Steps

**Step 1:** Select a (frozen) pretrained model and a layer/block to train on top of frozen features.

**Step 2:** Process images and cache the (frozen) features.

**Step 3:** Train a (lightweight) model on (augmented) frozen features.

$N$ : # image patches  
 $H$ : image height  
 $D$ : hidden size  
 $W$ : image width  
 $S$ : # classes  
 $C$ : # color channels  
 $L$ : # Transformer blocks

# Method Description

## Training on Top of Frozen Features in Three Steps

$N$ : # image patches

$H$ : image height

$D$ : hidden size

$W$ : image width

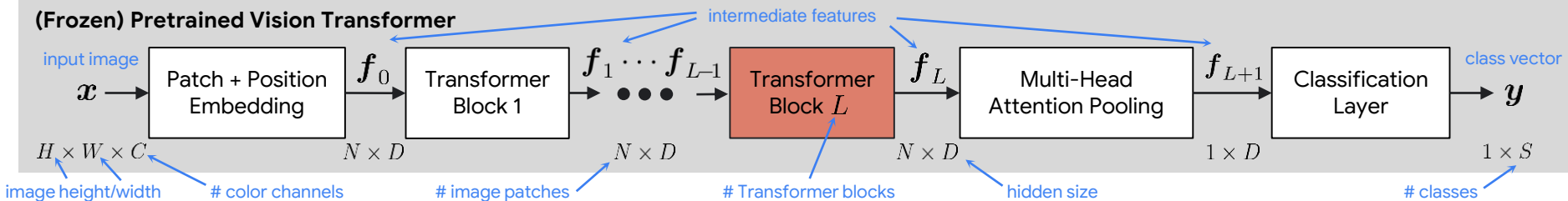
$S$ : # classes

$C$ : # color channels

$L$ : # Transformer blocks

**Step 1:** Select a (frozen) pretrained model and a **layer/block** to train on top of frozen features.

### (Frozen) Pretrained Vision Transformer



**Step 2:** Process images and cache the (frozen) features.

**Step 3:** Train a (lightweight) model on (augmented) frozen features.

# Method Description

## Training on Top of Frozen Features in Three Steps

$N$ : # image patches

$H$ : image height

$D$ : hidden size

$W$ : image width

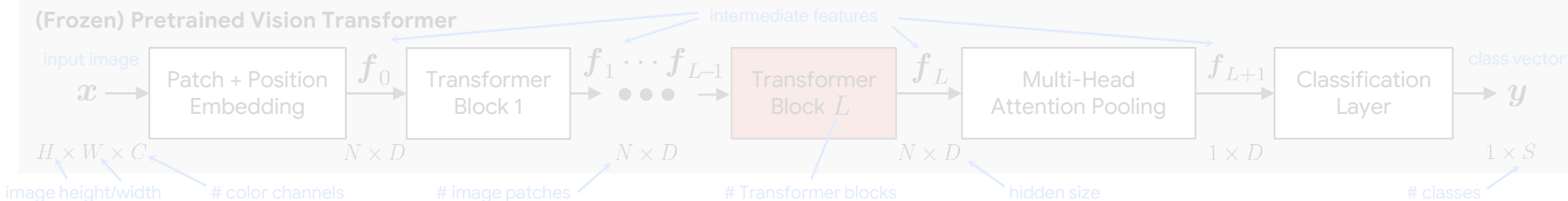
$S$ : # classes

$C$ : # color channels

$L$ : # Transformer blocks

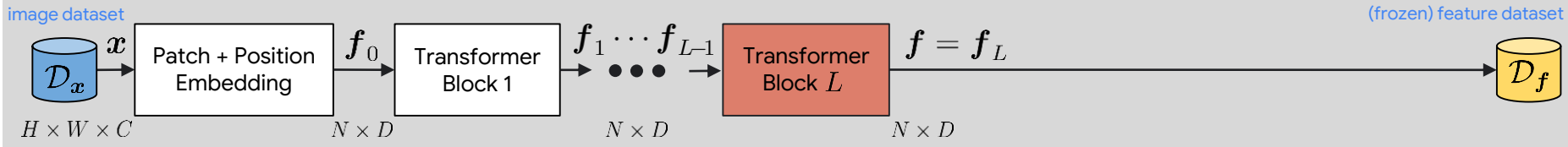
**Step 1:** Select a (frozen) pretrained model and a **layer/block** to train on top of frozen features.

### (Frozen) Pretrained Vision Transformer



**Step 2:** Process **images** and cache the **(frozen) features**.

### (Frozen) Pretrained Vision Transformer



**Step 3:** Train a (lightweight) model on (augmented) frozen features.

# Method Description

## Training on Top of Frozen Features in Three Steps

$N$ : # image patches

$H$ : image height

$D$ : hidden size

$W$ : image width

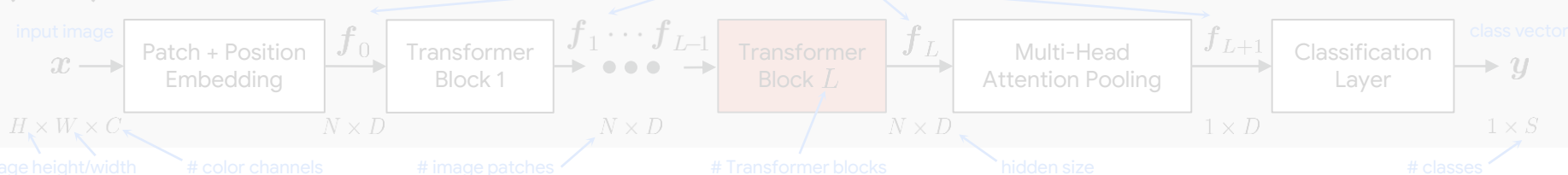
$S$ : # classes

$C$ : # color channels

$L$ : # Transformer blocks

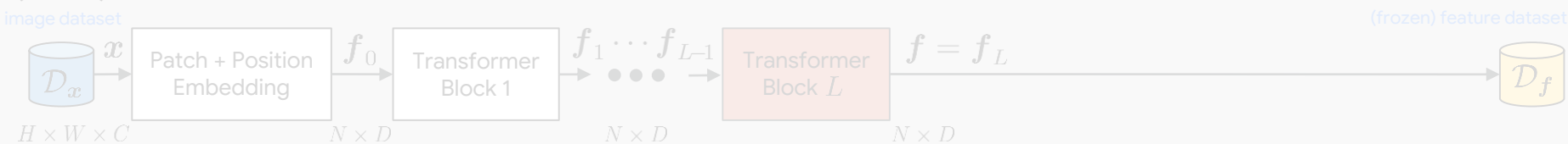
**Step 1:** Select a (frozen) pretrained model and a **layer/block** to train on top of frozen features.

### (Frozen) Pretrained Vision Transformer



**Step 2:** Process **images** and cache the **(frozen) features**.

### (Frozen) Pretrained Vision Transformer



**Step 3:** Train a **(lightweight) model** on **(augmented) frozen features**.

(frozen) feature dataset





# Method Description

## From Image Augmentations to Frozen Feature Augmentations (FroFAs)

**Step 0:** Feature reshaping from 2D ( $\mathbf{f}$ ) to 3D ( $\mathbf{f}^*$ ). - optional -

Reminder:  $\mathbf{x} \in \mathbb{I}^{H \times W \times C}$  ← # color channels ... with  $\mathbb{I} = [0, 1]$  or  $\mathbb{I} = \{0, 1, 2, \dots, 255\}$   
image height/width image value range

$\mathbf{f} = (\mathbf{f}_{n,d}) \in \mathbb{R}^{N \times D}$  ... reshaped to ...  $\mathbf{f}^* = (\mathbf{f}_{n_1, n_2, d}^*) \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D}$  ... with ...  $\mathbf{f}_d^* = \mathbf{f}_{::, d}^* \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times 1}$   
# image patches hidden size patch index channel index 2D patch indices 2D-ordered patches of a channel

# Method Description

## From Image Augmentations to Frozen Feature Augmentations (FroFAs)

Step 0: Feature reshaping from 2D ( $\mathbf{f}$ ) to 3D ( $\mathbf{f}^*$ ). - optional -

Reminder:  $\mathbf{x} \in \mathbb{I}^{H \times W \times C}$  ... with  $\mathbb{I} = [0, 1]$  or  $\mathbb{I} = \{0, 1, 2, \dots, 255\}$

*(Annotations:  $C$  is # color channels;  $H, W$  are image height/width;  $\mathbb{I}$  is image value range)*

$\mathbf{f} = (\mathbf{f}_{n,d}) \in \mathbb{R}^{N \times D}$  ... reshaped to ...  $\mathbf{f}^* = (\mathbf{f}_{n_1, n_2, d}^*) \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D}$  ... with ...  $\mathbf{f}_d^* = \mathbf{f}_{::,d}^* \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times 1}$

*(Annotations:  $N$  is # image patches;  $D$  is hidden size;  $n, d$  are patch index and channel index;  $n_1, n_2, d$  are 2D patch indices;  $d$  is 2D-ordered patches of a channel)*

Step 1: Address the channel dimensionality mismatch between  $\mathbf{x}$  and  $\mathbf{f}$  or  $\mathbf{f}^*$ .

# Method Description

## From Image Augmentations to Frozen Feature Augmentations (FroFAs)

Step 0: Feature reshaping from 2D ( $\mathbf{f}$ ) to 3D ( $\mathbf{f}^*$ ). - optional -

Reminder:  $\mathbf{x} \in \mathbb{I}^{H \times W \times C}$  ← # color channels ... with  $\mathbb{I} = [0, 1]$  or  $\mathbb{I} = \{0, 1, 2, \dots, 255\}$   
image height/width image value range

$\mathbf{f} = (\mathbf{f}_{n,d}) \in \mathbb{R}^{N \times D}$  ← # image patches hidden size ... reshaped to ...  $\mathbf{f}^* = (\mathbf{f}_{n_1, n_2, d}^*) \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D}$  ... with ...  $\mathbf{f}_d^* = \mathbf{f}_{:::,d}^* \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times \mathbb{I}}$   
patch index channel index 2D patch indices 2D-ordered patches of a channel

Step 1: Address the channel dimensionality mismatch between  $\mathbf{x}$  and  $\mathbf{f}$  or  $\mathbf{f}^*$ .

- Most image augmentations can be applied channel-wise, e.g., brightness adjustments.
- We re-use these augmentations in the feature space and ignore image augmentations relying on three color channels, e.g., color jitter.

# Method Description

## From Image Augmentations to Frozen Feature Augmentations (FroFAs)

Step 0: Feature reshaping from 2D ( $f$ ) to 3D ( $f^*$ ). - optional -

Reminder:  $\mathbf{x} \in \mathbb{I}^{H \times W \times C}$  ← # color channels ... with  $\mathbb{I} = [0, 1]$  or  $\mathbb{I} = \{0, 1, 2, \dots, 255\}$   
image height/width image value range

$\mathbf{f} = (f_{n,d}) \in \mathbb{R}^{N \times D}$  ← # image patches hidden size ... reshaped to ...  $\mathbf{f}^* = (f^*_{n_1, n_2, d}) \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D}$  ... with ...  $\mathbf{f}_d^* = f^*_{:::, d} \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times \mathbb{I}}$   
patch index channel index 2D patch indices 2D-ordered patches of a channel

Step 1: Address the channel dimensionality mismatch between  $\mathbf{x}$  and  $\mathbf{f}$  or  $\mathbf{f}^*$ .

- Most image augmentations can be applied channel-wise, e.g., brightness adjustments.
- We re-use these augmentations in the feature space and ignore image augmentations relying on three color channels, e.g., color jitter.

Step 2: Address the value range mismatch between  $\mathbf{x}$  and  $\mathbf{f}$  or  $\mathbf{f}^*$ .

# Method Description

## From Image Augmentations to Frozen Feature Augmentations (FroFAs)

**Step 0:** Feature reshaping from 2D ( $f$ ) to 3D ( $f^*$ ). - optional -

Reminder:  $x \in \mathbb{I}^{H \times W \times C}$  ← # color channels ... with  $\mathbb{I} = [0, 1]$  or  $\mathbb{I} = \{0, 1, 2, \dots, 255\}$   
image height/width image value range

$f = (f_{n,d}) \in \mathbb{R}^{N \times D}$  ... reshaped to ...  $f^* = (f^*_{n_1, n_2, d}) \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D}$  ... with ...  $f^*_d = f^*_{:, :, d} \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times \mathbb{I}}$   
# image patches hidden size patch index channel index 2D patch indices 2D-ordered patches of a channel

**Step 1:** Address the channel dimensionality mismatch between  $x$  and  $f$  or  $f^*$ .

- Most image augmentations can be applied channel-wise, e.g., brightness adjustments.
- We re-use these augmentations in the feature space and ignore image augmentations relying on three color channels, e.g., color jitter.

**Step 2:** Address the value range mismatch between  $x$  and  $f$  or  $f^*$ .

- We define feature-to-image/image-to feature transformations applied simultaneously on  $D_t$  channels.

$$t_{f \rightarrow x} : \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D_t} \rightarrow \mathbb{I}^{\sqrt{N} \times \sqrt{N} \times D_t} \quad t_{f \leftarrow x} : \mathbb{I}^{\sqrt{N} \times \sqrt{N} \times D_t} \rightarrow \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D_t}$$

We perform a linear scaling:

$$x_f = t_{f \rightarrow x}(f^*) = \frac{f^* - f_{\min}}{f_{\max} - f_{\min}}$$

# Method Description

## From Image Augmentations to Frozen Feature Augmentations (FroFAs)

**Step 0:** Feature reshaping from 2D ( $f$ ) to 3D ( $f^*$ ). - optional -

Reminder:  $x \in \mathbb{I}^{H \times W \times C}$  ← # color channels ... with  $\mathbb{I} = [0, 1]$  or  $\mathbb{I} = \{0, 1, 2, \dots, 255\}$   
 image height/width image value range

$f = (f_{n,d}) \in \mathbb{R}^{N \times D}$  ... reshaped to ...  $f^* = (f^*_{n_1, n_2, d}) \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D}$  ... with ...  $f^*_d = f^*_{:, :, d} \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times \mathbb{I}}$   
 # image patches hidden size patch index channel index 2D patch indices 2D-ordered patches of a channel

**Step 1:** Address the channel dimensionality mismatch between  $x$  and  $f$  or  $f^*$ .

- Most image augmentations can be applied channel-wise, e.g., brightness adjustments.
- We re-use these augmentations in the feature space and ignore image augmentations relying on three color channels, e.g., color jitter.

**Step 2:** Address the value range mismatch between  $x$  and  $f$  or  $f^*$ .

- We define feature-to-image/image-to feature transformations applied simultaneously on  $D_t$  channels.

$$t_{f \rightarrow x} : \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D_t} \rightarrow \mathbb{I}^{\sqrt{N} \times \sqrt{N} \times D_t} \quad t_{f \leftarrow x} : \mathbb{I}^{\sqrt{N} \times \sqrt{N} \times D_t} \rightarrow \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D_t}$$

We perform a linear scaling:

$$x_f = t_{f \rightarrow x}(f^*) = \frac{f^* - f_{\min}}{f_{\max} - f_{\min}}$$

**Step 3:** Putting all together → FroFA ( $a_f$ ):

$$a_f = t_{f \leftarrow x} \circ a_x \circ t_{f \rightarrow x} \quad \dots \text{with} \quad a_x : \mathbb{I}^{\sqrt{N} \times \sqrt{N} \times D_a} \rightarrow \mathbb{I}^{\sqrt{N} \times \sqrt{N} \times D_a}$$

image augmentation

# Method Description

## From Image Augmentations to Frozen Feature Augmentations (FroFAs)

**Step 0:** Feature reshaping from 2D ( $f$ ) to 3D ( $f^*$ ). - optional -

Reminder:  $x \in \mathbb{I}^{H \times W \times C}$  ← # color channels ... with  $\mathbb{I} = [0, 1]$  or  $\mathbb{I} = \{0, 1, 2, \dots, 255\}$   
 image height/width image value range

$f = (f_{n,d}) \in \mathbb{R}^{N \times D}$  ... reshaped to ...  $f^* = (f^*_{n_1, n_2, d}) \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D}$  ... with ...  $f^*_d = f^*_{:, :, d} \in \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times \mathbb{I}}$   
 # image patches hidden size patch index channel index 2D patch indices 2D-ordered patches of a channel

**Step 1:** Address the channel dimensionality mismatch between  $x$  and  $f$  or  $f^*$ .

- Most image augmentations can be applied channel-wise, e.g., brightness adjustments.
- We re-use these augmentations in the feature space and ignore image augmentations relying on three color channels, e.g., color jitter.

**Step 2:** Address the value range mismatch between  $x$  and  $f$  or  $f^*$ .

- We define feature-to-image/image-to feature transformations applied simultaneously on  $D_t$  channels.

$$t_{f \rightarrow x} : \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D_t} \rightarrow \mathbb{I}^{\sqrt{N} \times \sqrt{N} \times D_t} \quad t_{f \leftarrow x} : \mathbb{I}^{\sqrt{N} \times \sqrt{N} \times D_t} \rightarrow \mathbb{R}^{\sqrt{N} \times \sqrt{N} \times D_t}$$

We perform a linear scaling:

$$x_f = t_{f \rightarrow x}(f^*) = \frac{f^* - f_{\min}}{f_{\max} - f_{\min}}$$

**Step 3:** Putting all together → FroFA ( $a_f$ ):

$$a_f = t_{f \leftarrow x} \circ a_x \circ t_{f \rightarrow x} \quad \dots \text{with} \quad a_x : \mathbb{I}^{\sqrt{N} \times \sqrt{N} \times D_a} \rightarrow \mathbb{I}^{\sqrt{N} \times \sqrt{N} \times D_a}$$

image augmentation

By default, we have the same mapping and randomness across channels.

We tested channel-wise interpretations:

- Channel-wise randomness
- Channel-wise mapping/randomness

# Experimental Setup

Network Architectures, Baseline Models, Pretraining/Few-Shot Transfer Datasets, Data Augmentation

Network architectures:

ViT Model	# Params (M)
Ti/16	5.5
B/16	86.0
L/16	303.0

Pure **image classification** pretraining: **Vision-language** pretraining:

Dataset	# Images
JFT-3B	~ 3 billion
ImageNet-21k	~ 14 million

Dataset	# Images-text pairs
WebLI	~ 31 billion

Few-shot learning (1-, 5-, 10-, 25-shot):

Dataset	# Classes
ImageNet1k	1000
CIFAR10	10
CIFAR100	100
DMLab	6
DTD	47
Resisc45	45
SUN397	397
SVHN	10



# Experimental Setup

Network Architectures, Baseline Models, Pretraining/Few-Shot Transfer Datasets, Data Augmentation

Network architectures:

ViT Model	# Params (M)
Ti/16	5.5
B/16	86.0
L/16	303.0

Pure **image classification** pretraining: **Vision-language** pretraining:

Dataset	# Images
JFT-3B	~ 3 billion
ImageNet-21k	~ 14 million

Dataset	# Images-text pairs
WebLI	~ 31 billion

Few-shot learning (1-, 5-, 10-, 25-shot):

Dataset	# Classes
ImageNet1k	1000
CIFAR10	10
CIFAR100	100
DMLab	6
DTD	47
Resisc45	45
SUN397	397
SVHN	10

- Other**
- JPEG
  - mixup

- Crop & Drop**
- crop
  - resized crop
  - inception crop
  - patch dropout

Baseline models:

Baseline	Input size	Regularization
MAP <sup>wd</sup>	$N \times D$	weight decay
Linear probe	$1 \times D$	L2 regularization

our lightweight model, but without FroFA

penultimate output of the pretrained model

- Stylistic**
- brightness
  - contrast
  - equalize
  - invert
  - posterize
  - sharpness
  - solarize

Data augmentation

- Geometric**
- rotate
  - shear-x,
  - shear-y,
  - translate-x,
  - translate-y

# Experimental Setup

Network Architectures, Baseline Models, Pretraining/Few-Shot Transfer Datasets, Data Augmentation

Network architectures:

ViT Model	# Params (M)
Ti/16	5.5
B/16	86.0
L/16	303.0

Pure **image classification** pretraining: **Vision-language** pretraining:

Dataset	# Images
JFT-3B	~ 3 billion
ImageNet-21k	~ 14 million

Dataset	# Images-text pairs
WebLI	~ 31 billion

Few-shot learning (1-, 5-, 10-, 25-shot):

Dataset	# Classes
ImageNet1k	1000
CIFAR10	10
CIFAR100	100
DMLab	6
DTD	47
Resisc45	45
SUN397	397
SVHN	10

Baseline models:

Baseline	Input size	Regularization
MAP <sup>wd</sup>	$N \times D$	weight decay
Linear probe	$1 \times D$	L2 regularization

our lightweight model, but without FroFA

penultimate output of the pretrained model

- Other**
- JPEG
  - mixup

- Crop & Drop**
- crop
  - resized crop
  - inception crop
  - patch dropout

- Stylistic**
- brightness
  - contrast
  - equalize
  - invert
  - posterize
  - sharpness
  - solarize

Data augmentation

- Geometric**
- rotate
  - shear-x,
  - shear-y,
  - translate-x,
  - translate-y

**TL;DR**

We test **three** ViT models, **three** pretraining datasets, **eight** few-shot datasets, **eighteen** data augmentations ...  
... and compare to **two** baselines.

# Experimental Setup

Network Architectures, Baseline Models, Pretraining/Few-Shot Transfer Datasets, Data Augmentation

Network architectures:

ViT Model	# Params (M)
Ti/16	5.5
B/16	86.0
<b>L/16</b>	303.0

Pure image classification pretraining: Vision-language pretraining:

Dataset	# Images
<b>JFT-3B</b>	~ 3 billion
ImageNet-21k	~ 14 million

Dataset	# Images-text pairs
<b>WebLI</b>	~ 31 billion

Few-shot learning (1-, 5-, 10-, 25-shot):

Dataset	# Classes
ImageNet1k	1000
<b>CIFAR10</b>	10
<b>CIFAR100</b>	100
<b>DMLab</b>	6
<b>DTD</b>	47
<b>Resisc45</b>	45
<b>SUN397</b>	397
<b>SVHN</b>	10

Baseline models:

Baseline	Input size	Regularization
MAP <sup>wd</sup>	$N \times D$	weight decay
Linear probe	$1 \times D$	L2 regularization

our lightweight model, but without FroFA

penultimate output of the pretrained model

- Other
- JPEG
  - mixup

- Crop & Drop
- crop
  - resized crop
  - inception crop
  - patch dropout

- Stylistic
- **brightness**
  - contrast
  - equalize
  - invert
  - posterize
  - sharpness
  - solarize

Data augmentation

- Geometric
- rotate
  - shear-x,
  - shear-y,
  - translate-x,
  - translate-y

TL;DR

We test **three** ViT models, **three** pretraining datasets, **eight** few-shot datasets, **eighteen** data augmentations ...  
... and compare to **two** baselines.

In the following, we focus on a subset of our results.

# Experimental Results (Highlights)

L/16 Pretrained on **JFT-3B** or **WebLI** and Transferred to Seven Few-Shot Datasets Using a Channel Variant of Brightness FroFA

We compute the average gains across **7 few-shot datasets (SUN397, ...)**, excluding ImageNet1k.

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25
	<b>MAP<sup>wd</sup></b>				
<b>JFT-3B</b>	<b>Linear probe</b>				
	<b>MAP<sup>wd</sup> + FroFA (ours)</b>				

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25
	<b>MAP<sup>wd</sup></b>				
<b>WebLI + SigLIP</b>	<b>Linear probe</b>				
	<b>MAP<sup>wd</sup> + FroFA (ours)</b>				

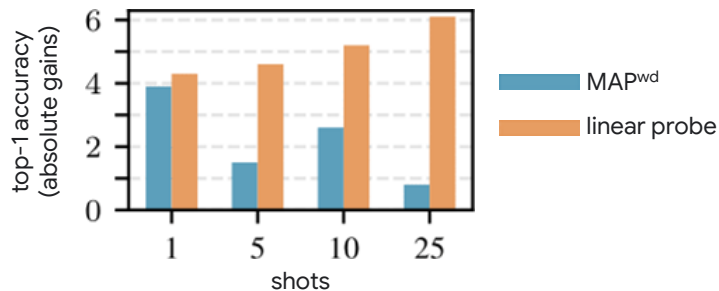
# Experimental Results (Highlights)

L/16 Pretrained on **JFT-3B** or **WebLI** and Transferred to Seven Few-Shot Datasets Using a Channel Variant of Brightness FroFA

We compute the average gains across **7 few-shot datasets (SUN397, ...)**, excluding ImageNet1k.

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25
JFT-3B	MAP <sup>wd</sup>	49.5	65.8	68.3	74.1
	Linear probe	49.1	62.7	65.7	68.8
	MAP <sup>wd</sup> + FroFA (ours)	<b>53.4</b>	<b>67.3</b>	<b>70.9</b>	<b>74.9</b>

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25
WebLI + SigLIP	MAP <sup>wd</sup>				
	Linear probe				
	MAP <sup>wd</sup> + FroFA (ours)				



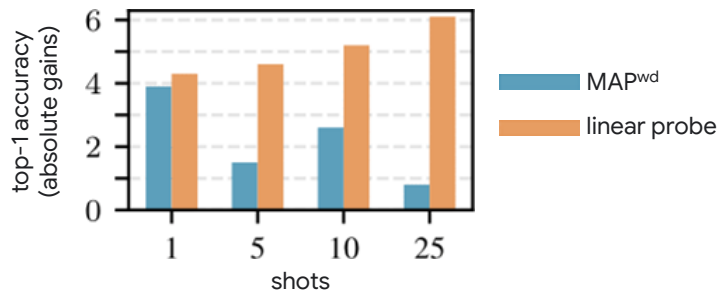
# Experimental Results (Highlights)

L/16 Pretrained on **JFT-3B** or **WebLI** and Transferred to Seven Few-Shot Datasets Using a Channel Variant of Brightness FroFA

We compute the average gains across **7 few-shot datasets (SUN397, ...)**, excluding ImageNet1k.

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25
JFT-3B	MAP <sup>wd</sup>	49.5	65.8	68.3	74.1
	Linear probe	49.1	62.7	65.7	68.8
	MAP <sup>wd</sup> + FroFA (ours)	<b>53.4</b>	<b>67.3</b>	<b>70.9</b>	<b>74.9</b>

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25
WebLI + SigLIP	MAP <sup>wd</sup>				
	Linear probe				
	MAP <sup>wd</sup> + FroFA (ours)				



JFT-3B pretraining:

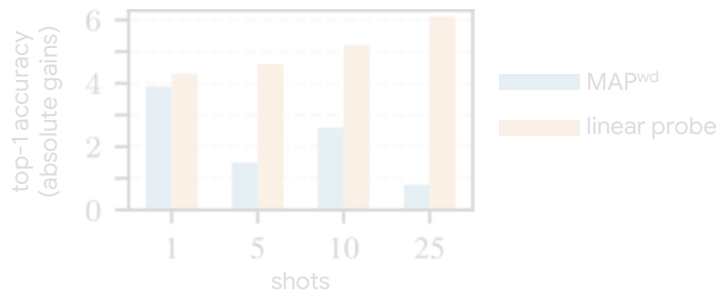
- On average, we **improve** upon both baselines (MAP<sup>wd</sup>, linear probe).
- The **gains** for MAP<sup>wd</sup> **diminish** with **higher shots**, but for **linear probe** **increase** with **higher shots**.

# Experimental Results (Highlights)

L/16 Pretrained on **JFT-3B** or **WebLI** and Transferred to Seven Few-Shot Datasets Using a Channel Variant of Brightness FroFA

We compute the average gains across **7 few-shot datasets (SUN397, ...)**, excluding ImageNet1k.

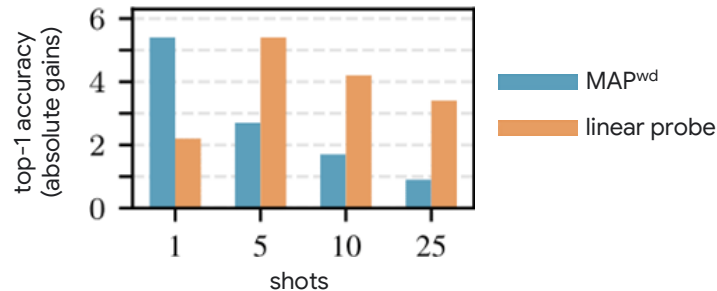
top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25
JFT-3B	MAP <sup>wd</sup>	49.5	65.8	68.3	74.1
	Linear probe	49.1	62.7	65.7	68.8
	MAP <sup>wd</sup> + FroFA (ours)	<b>53.4</b>	<b>67.3</b>	<b>70.9</b>	<b>74.9</b>



JFT-3B pretraining:

- On average, we **improve** upon both baselines (MAP<sup>wd</sup>, linear probe).
- The **gains** for MAP<sup>wd</sup> **diminish** with **higher shots**, but for **linear probe** **increase** with **higher shots**.

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25
WebLI + SigLIP	MAP <sup>wd</sup>	45.9	67.7	71.8	75.1
	Linear probe	49.1	65.0	69.3	72.6
	MAP <sup>wd</sup> + FroFA (ours)	<b>51.3</b>	<b>70.4</b>	<b>73.5</b>	<b>76.0</b>

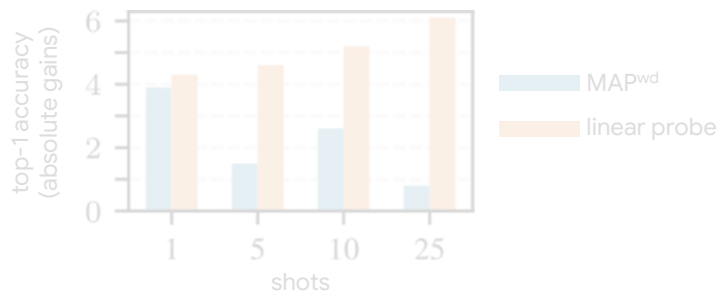


# Experimental Results (Highlights)

L/16 Pretrained on **JFT-3B** or **WebLI** and Transferred to Seven Few-Shot Datasets Using a Channel Variant of Brightness FroFA

We compute the average gains across **7 few-shot datasets (SUN397, ...)**, excluding ImageNet1k.

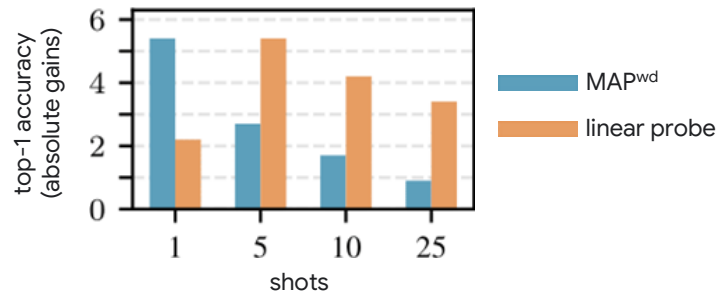
top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25
JFT-3B	MAP <sup>wd</sup>	49.5	65.8	68.3	74.1
	Linear probe	49.1	62.7	65.7	68.8
	MAP <sup>wd</sup> + FroFA (ours)	<b>53.4</b>	<b>67.3</b>	<b>70.9</b>	<b>74.9</b>



JFT-3B pretraining:

- On average, we **improve** upon both baselines (MAP<sup>wd</sup>, linear probe).
- The **gains** for MAP<sup>wd</sup> **diminish** with **higher shots**, but for **linear probe** **increase** with **higher shots**.

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25
WebLI + SigLIP	MAP <sup>wd</sup>	45.9	67.7	71.8	75.1
	Linear probe	49.1	65.0	69.3	72.6
	MAP <sup>wd</sup> + FroFA (ours)	<b>51.3</b>	<b>70.4</b>	<b>73.5</b>	<b>76.0</b>



WebLI vision-language pretraining (based on SigLIP [Zhai+, 2023]):

[Zhai et al., „Sigmoid Loss for Language-Image Pretraining“, ICCV, October 2023]

- On average, we **improve** upon both baselines (MAP<sup>wd</sup>, linear probe).
- The **gains** for both MAP<sup>wd</sup> and **linear probe** **diminish** with **higher shots**.



# Experimental Results (Highlights)

L/16 Pretrained on **JFT-3B** or **WebLI** and Transferred to Seven Few-Shot Datasets Using a Channel Variant of Brightness FroFA

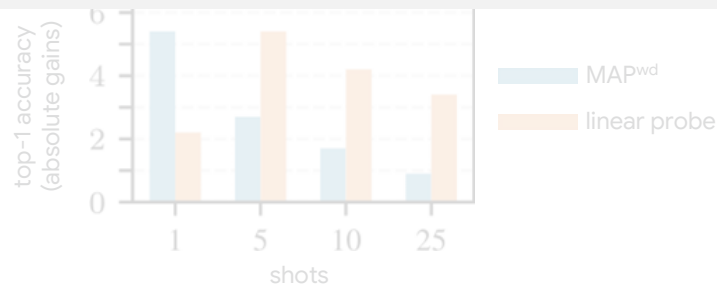
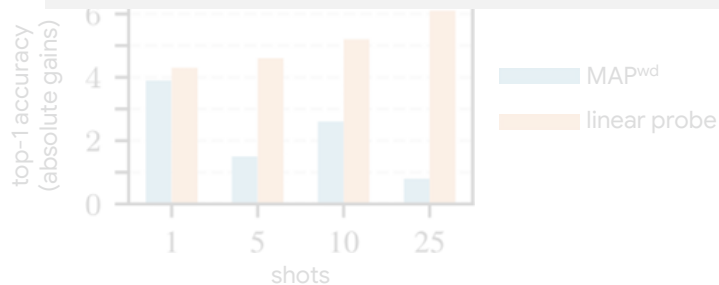
We compute the average gains across 7 few-shot datasets (SUN397, ...), excluding ImageNet1k.

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25

Initial Research question:

Can we effectively combine image augmentations and frozen features to train a lightweight model?



JFT-3B pretraining:

- On average, we **improve** upon both baselines (**MAP<sup>wd</sup>**, **linear probe**).
- The **gains** for **MAP<sup>wd</sup>** **diminish** with **higher shots**, but for **linear probe** **increase** with **higher shots**.

WebLI vision-language pretraining (based on SigLIP [Zhai+, 2023]):

[Zhai et al., „Sigmoid Loss for Language-Image Pretraining“, ICCV, October 2023]

- On average, we **improve** upon both baselines (**MAP<sup>wd</sup>**, **linear probe**).
- The **gains** for both **MAP<sup>wd</sup>** and **linear probe** **diminish** with **higher shots**.

# Experimental Results (Highlights)

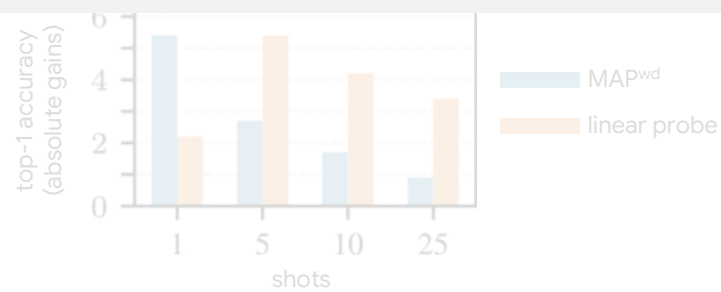
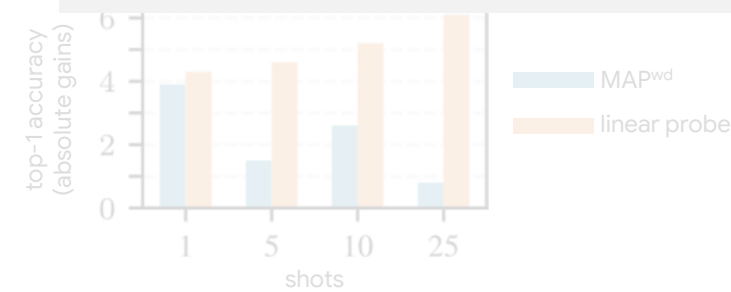
L/16 Pretrained on **JFT-3B** or **WebLI** and Transferred to Seven Few-Shot Datasets Using a Channel Variant of Brightness FroFA

We compute the average gains across 7 few-shot datasets (SUN397, ...), excluding ImageNet1k.

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25

top-1 accuracy [%]		Shots			
Pretraining scheme	Method	1	5	10	25

JFT	MAP <sup>wd</sup>	40.5	45.9	49.2	74.1	Initial Research question: Can we effectively combine image augmentations and frozen features to train a lightweight model?	45.9	47.7	71.9	75.1
	linear probe					Answer:				
						<b>Yes! Our method transfers across pretraining and few-shot transfer setups!</b>				
										72.6
										76.0



JFT-3B pretraining:

- On average, we **improve** upon both baselines (**MAP<sup>wd</sup>**, **linear probe**).
- The **gains** for **MAP<sup>wd</sup>** **diminish** with **higher shots**, but for **linear probe** **increase** with **higher shots**.

WebLI vision-language pretraining (based on SigLIP [Zhai+, 2023]):

- [Zhai et al., „Sigmoid Loss for Language-Image Pretraining“, ICCV, October 2023]
- On average, we **improve** upon both baselines (**MAP<sup>wd</sup>**, **linear probe**).
  - The **gains** for both **MAP<sup>wd</sup>** and **linear probe** **diminish** with **higher shots**.

# Conclusion

We investigated 18 frozen feature augmentations (FroFAs) for few-shot transfer learning for image classification.

We performed ablations along three axes: model size, pretraining, and transfer few-shot dataset.

Our main takeaways:

(a) Shown in this talk:

- Our best FroFA **excels** on **smaller few-shot datasets**.
- Our best FroFA **transfers** to **vision-language pretrained models**.

# Conclusion

We investigated 18 frozen feature augmentations (FroFAs) for few-shot transfer learning for image classification.

We performed ablations along three axes: model size, pretraining, and transfer few-shot dataset.

Our main takeaways:

(a) Shown in this talk:

- Our best FroFA **excels** on **smaller few-shot datasets**.
- Our best FroFA **transfers** to **vision-language pretrained models**.

(b) Not shown in this talk:

- Training with **stylistic FroFAs** give **large improvements** upon a representative baseline across all shots.
- **Channel variants** can yield to **further improvements**.
- **Sequential protocols** seem **promising** for future works, given our simple proof of concept.

For more details, checkout our paper or let's have a chat at the conference!



# Thank you.

Link to our paper:



Google DeepMind



**Andreas Bär**  
*TU Braunschweig*



**Neil Houlsby**  
*Google DeepMind*



**Mostafa Dehghani**  
*Google DeepMind*



**Manoj Kumar**  
*Google DeepMind*