

Cornell Bowers C-IS

College of Computing and Information Science

Deep Learning

Week 05: ViTs and CLIP

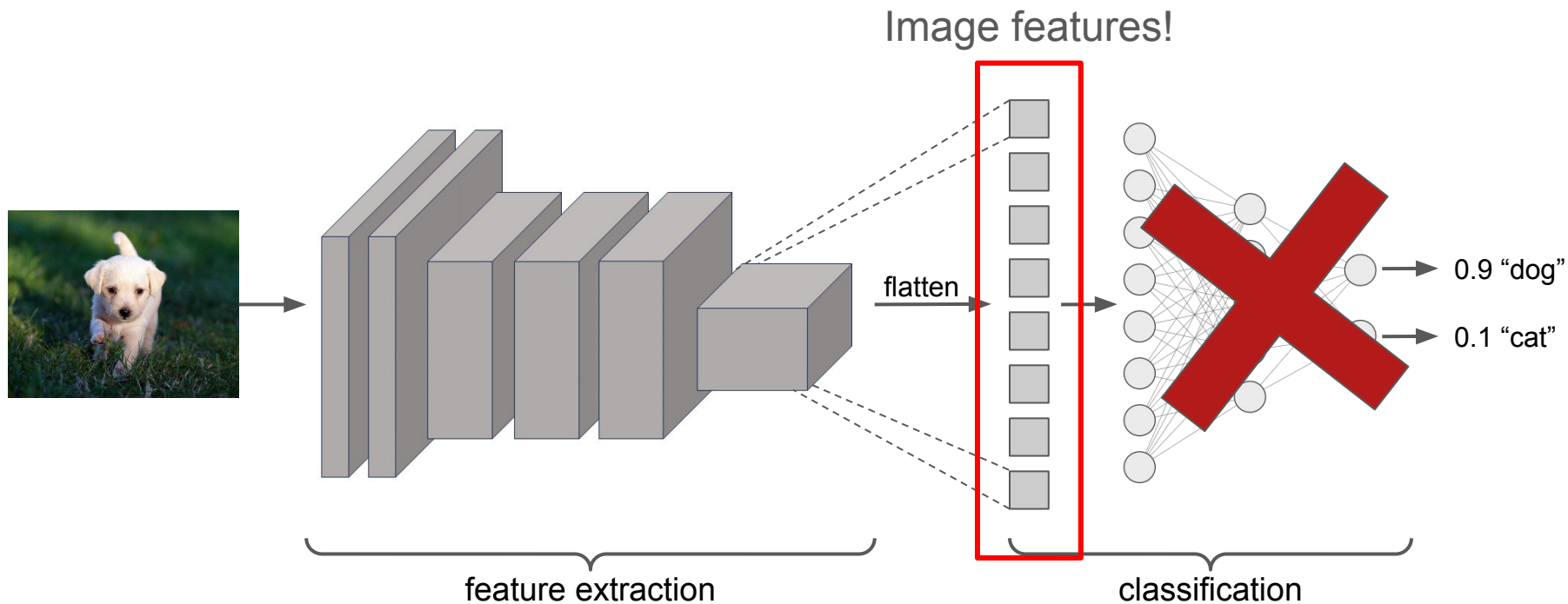
Thanks to

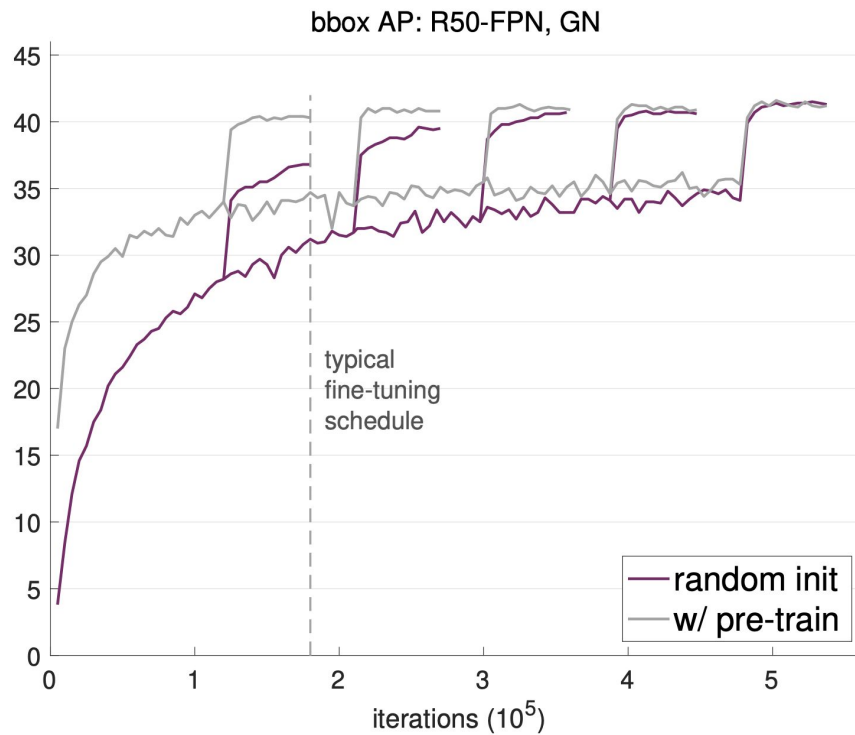
Varsha Kishore
Justin Lovelace
Anissa Dallmann
Dylan Van Bramer

Logistics

- **HW2+P2** is due today!
- **HW3+P3** will be released today. due 3/13/25 (**only 1 week!**)
- **Quiz 3** Paper released. Paper: *CLIP*
- **Final Project** signups released today. Look out for Ed announcement.

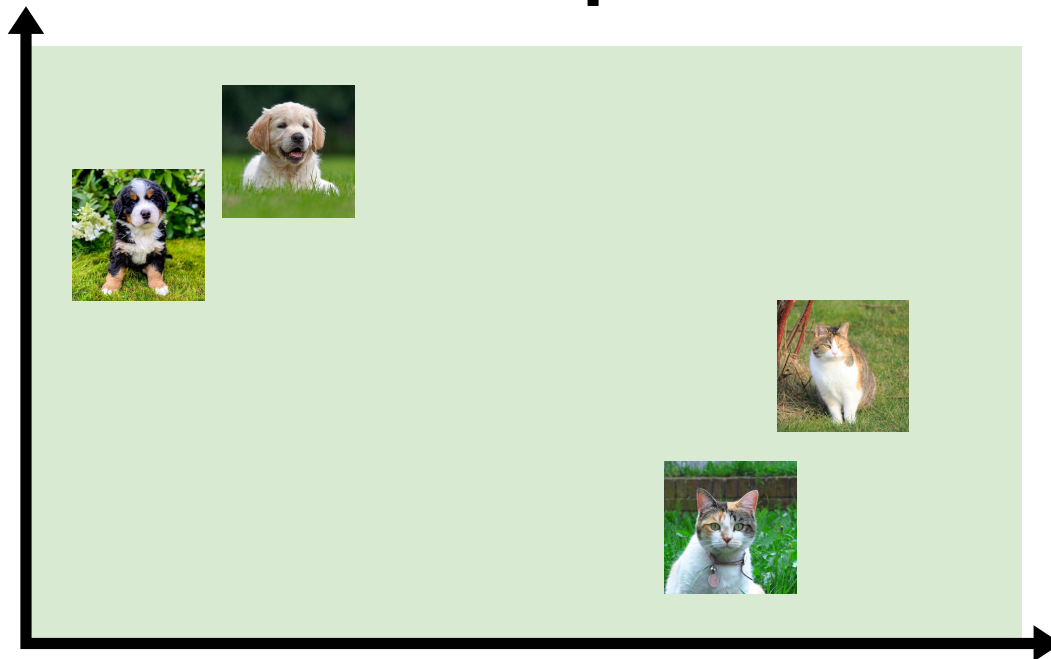
Image Classification





Use self-supervised learning to learn embeddings for images

Vector Space



Discussion: Comparison of Loss Functions

Triplet:

$$l = \max(0, \text{sim}(\mathbf{x}, \mathbf{x}^-) - \text{sim}(\mathbf{x}, \mathbf{x}^+) + c)$$

SimCLR loss:

$$l = -\log \left(\frac{\exp(\text{sim}(\mathbf{x}, \mathbf{x}^+)/\tau)}{\exp(\text{sim}(\mathbf{x}, \mathbf{x}^+)/\tau) + \exp(\text{sim}(\mathbf{x}, \mathbf{x}^-)/\tau)} \right)$$

Algorithm 1 SimCLR's main learning algorithm.

input: batch size N , constant τ , structure of f, g, \mathcal{T} .

for sampled minibatch $\{\mathbf{x}_k\}_{k=1}^N$ **do**

for all $k \in \{1, \dots, N\}$ **do**

 draw two augmentation functions $t \sim \mathcal{T}, t' \sim \mathcal{T}$

 # the first augmentation

$\tilde{\mathbf{x}}_{2k-1} = t(\mathbf{x}_k)$

$\mathbf{h}_{2k-1} = f(\tilde{\mathbf{x}}_{2k-1})$

 # representation

$\mathbf{z}_{2k-1} = g(\mathbf{h}_{2k-1})$

 # projection

 # the second augmentation

$\tilde{\mathbf{x}}_{2k} = t'(\mathbf{x}_k)$

$\mathbf{h}_{2k} = f(\tilde{\mathbf{x}}_{2k})$

 # representation

$\mathbf{z}_{2k} = g(\mathbf{h}_{2k})$

 # projection

end for

for all $i \in \{1, \dots, 2N\}$ and $j \in \{1, \dots, 2N\}$ **do**

$s_{i,j} = \mathbf{z}_i^\top \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$ # pairwise similarity

end for

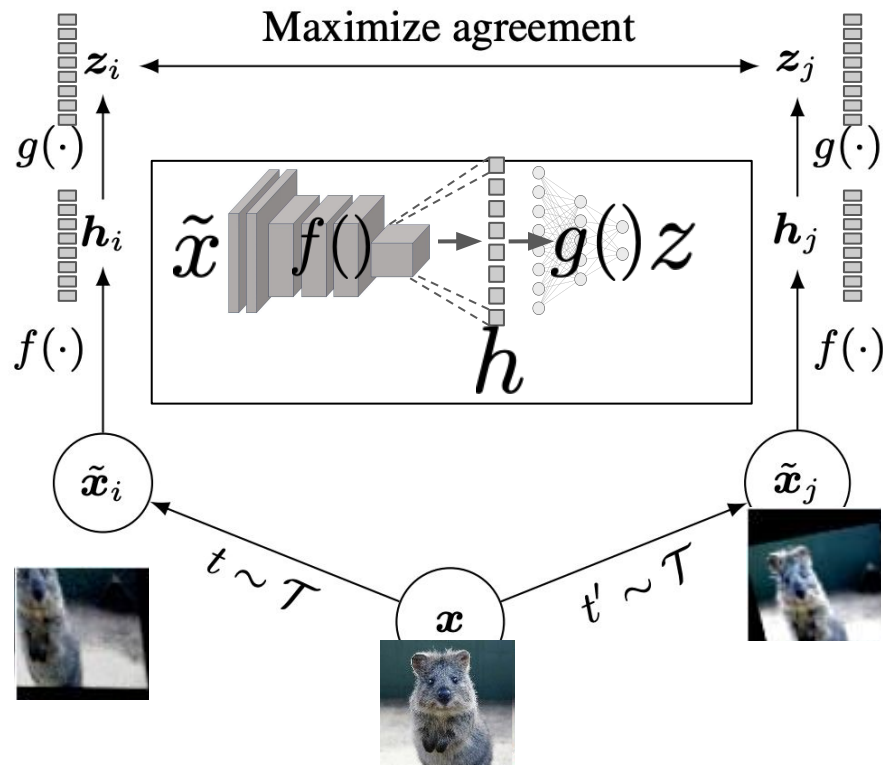
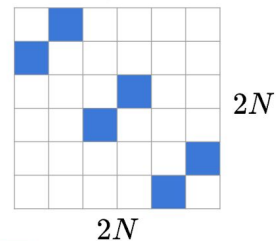
define $\ell(i, j)$ **as** $\ell(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(s_{i,k}/\tau)}$

$\mathcal{L} = \frac{1}{2N} \sum_{k=1}^N [\ell(2k-1, 2k) + \ell(2k, 2k-1)]$

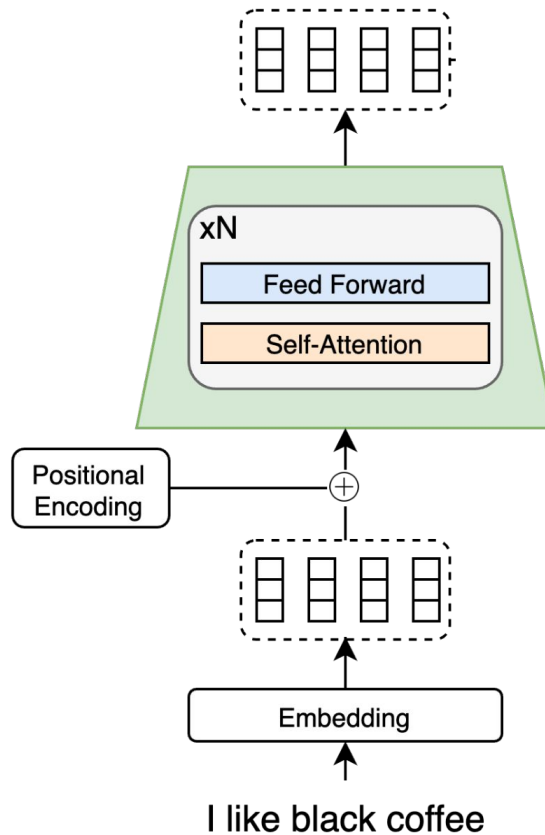
 update networks f and g to minimize \mathcal{L}

end for

return encoder network $f(\cdot)$, and throw away $g(\cdot)$



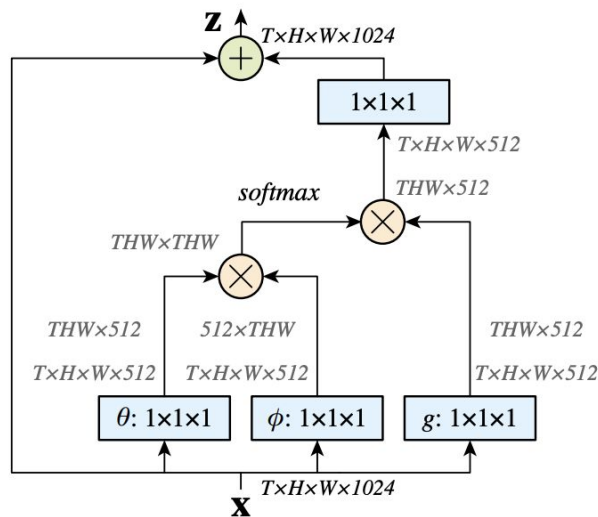
Remember Transformers?



How can we use Transformers on Images?

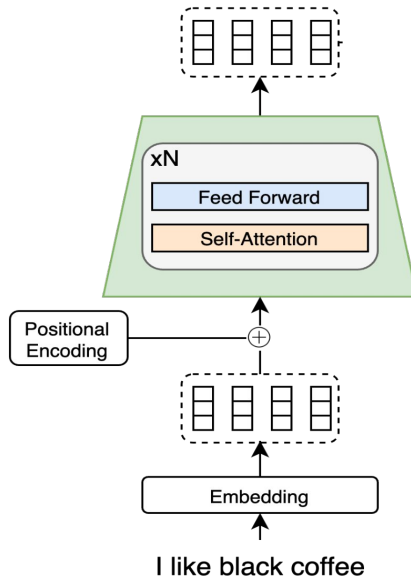
How to use Attention for Vision Tasks?

Attempt #1: Add attention to existing CNNs



How to use Attention for Vision Tasks?

Attempt #2: Adapt standard transformers to image data



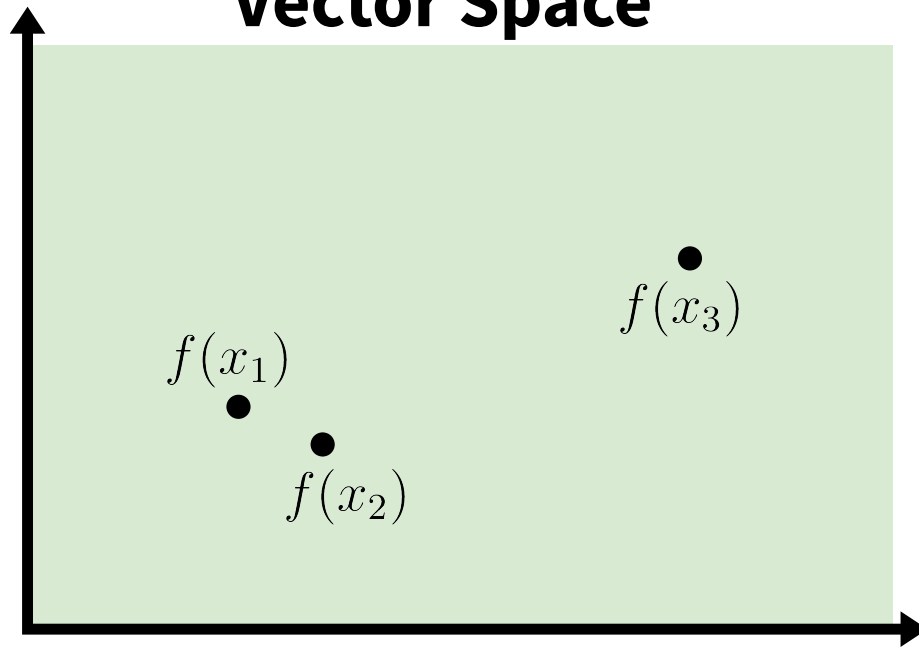
$f(x)$ = word embedding

“machine” x_1

“learning” x_2

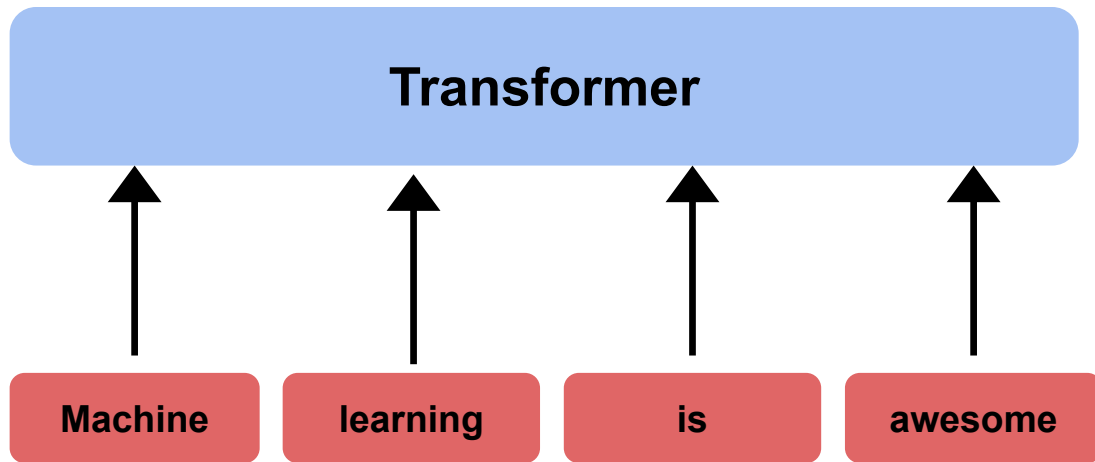
“awesome” x_3

Vector Space

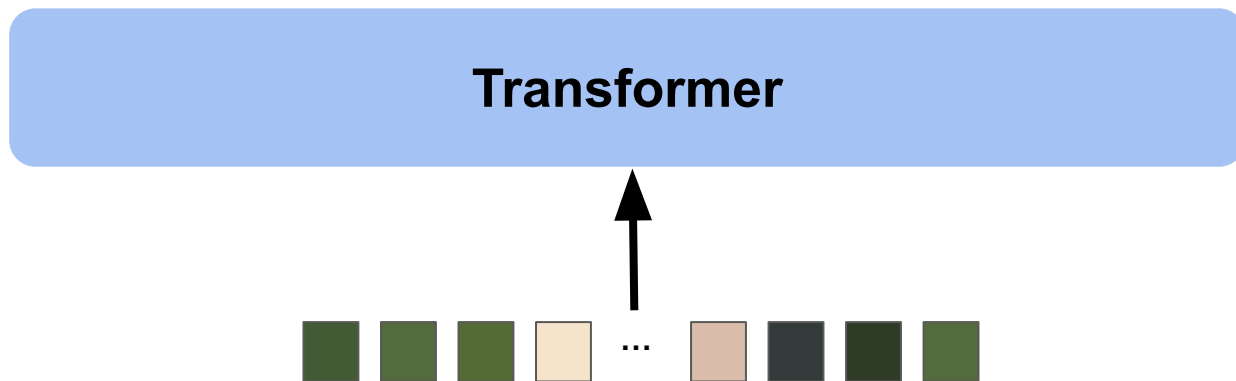


Can we extend this
idea to images?

Machine learning
is awesome



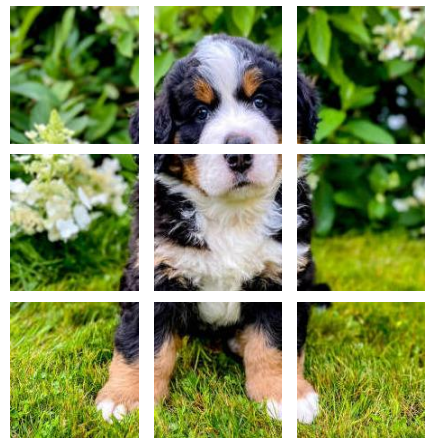
Idea1: Use pixels as input tokens



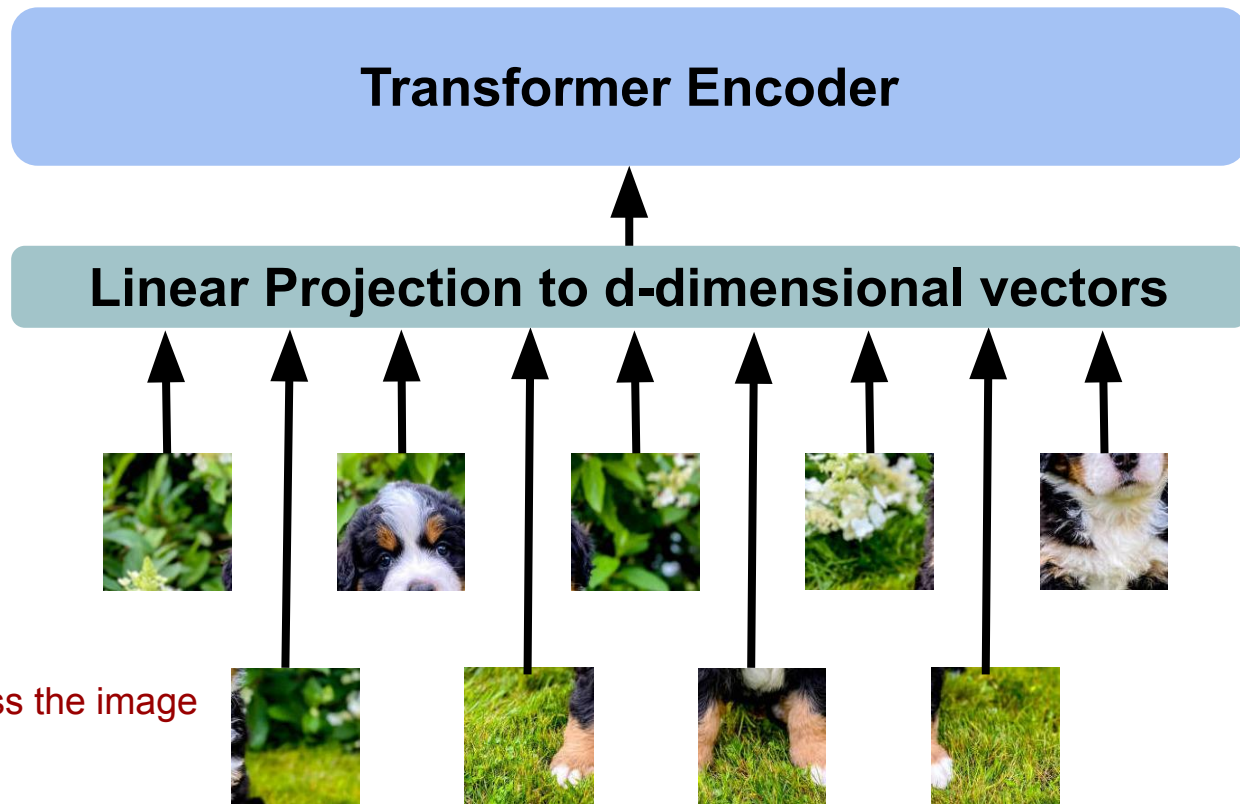
Discuss:

- How big would the attention matrix be for an $R \times R$ image?
- How much memory do we need for a 512×512 image, with 16 attention heads, and 48 layers - for attention alone?
- What if the image was 32×32 pixels?

Idea2: use image patches as input tokens



CNNs can be used to preprocess the image before converting it to patches!





x_1

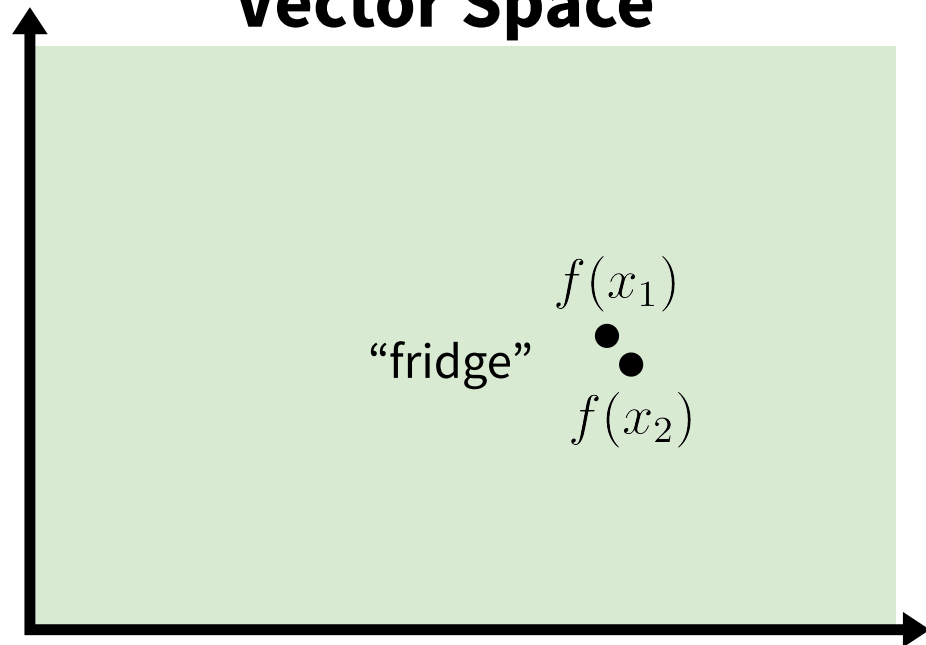


x_2

At a high level, these are both
just fridges

$$f(x) = \text{convolution}$$

Vector Space





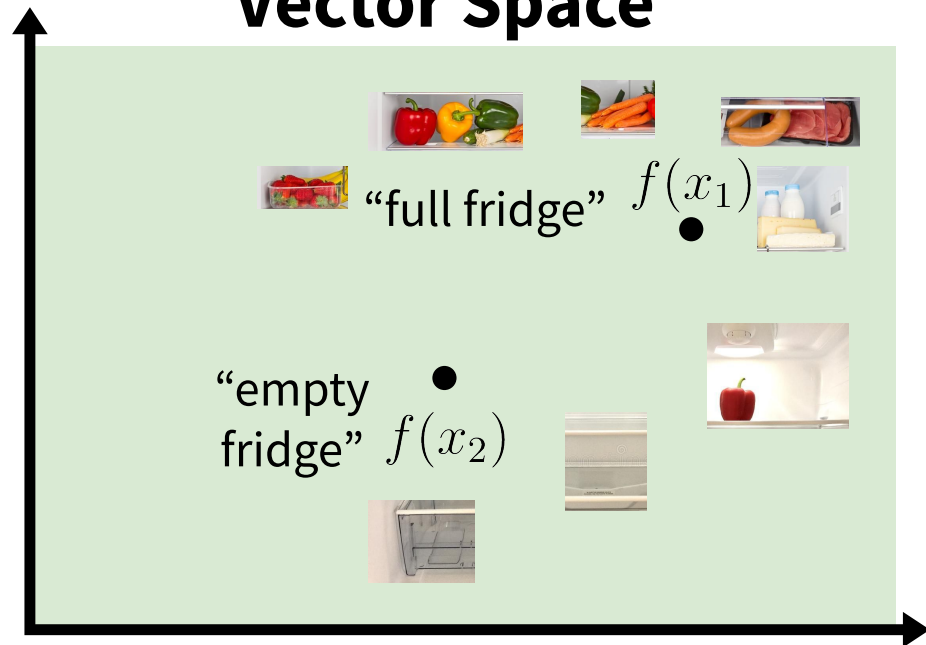
x_1

x_2

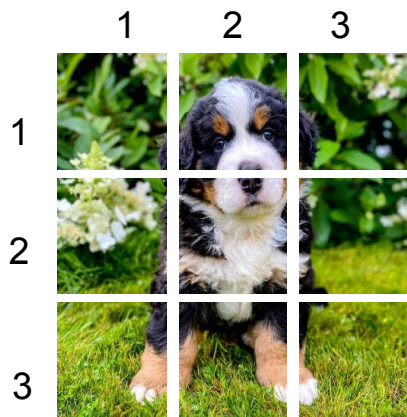
Use image patches like words in a sentence!

$f(x) =$ Vision Transformers

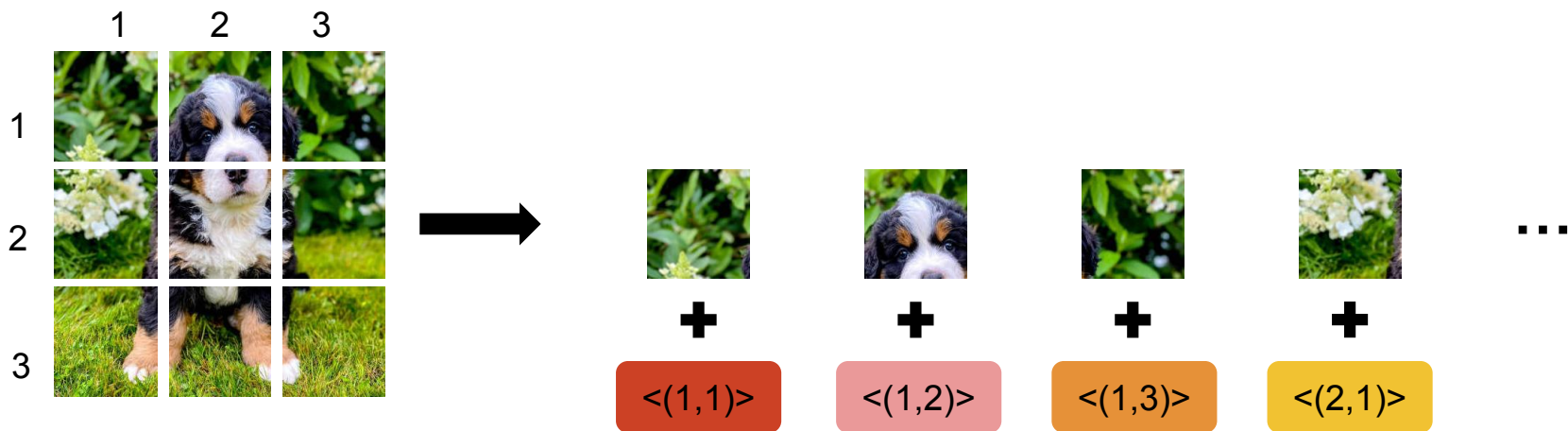
Vector Space



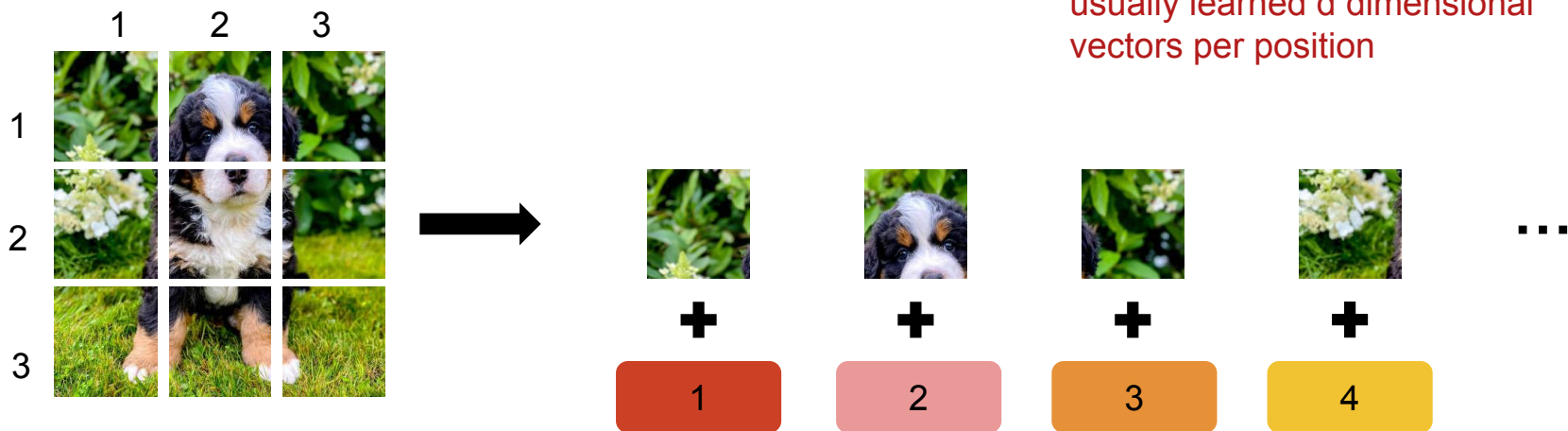
Adding positional embeddings



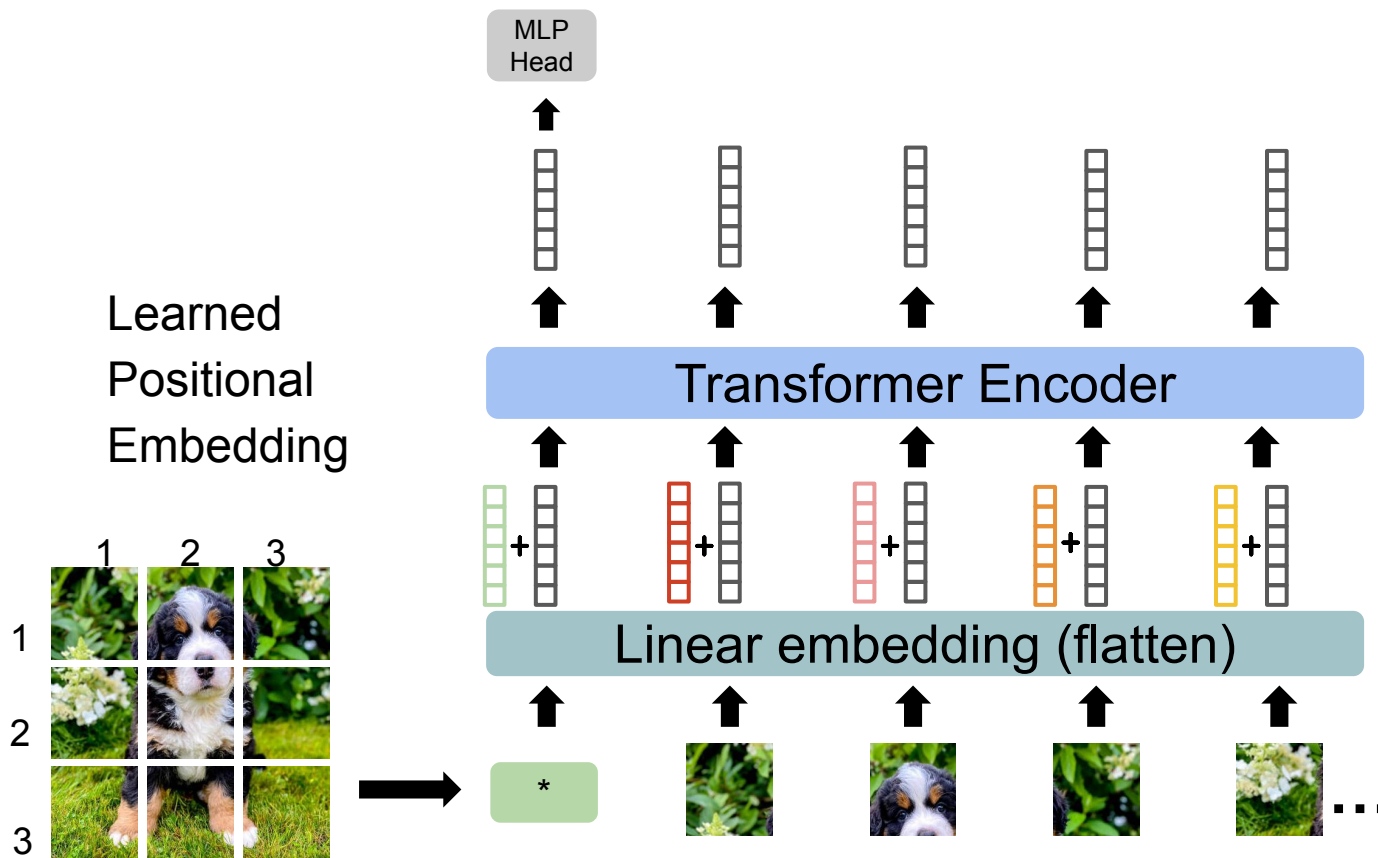
Adding positional embeddings



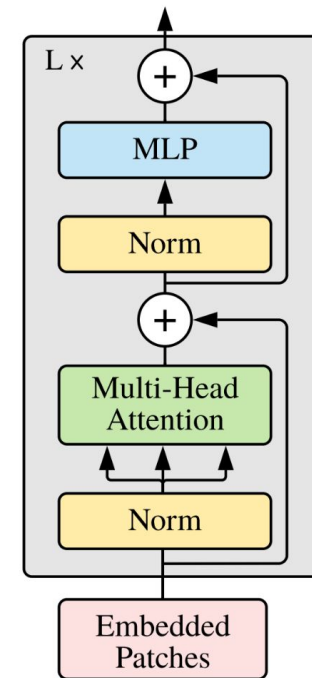
Adding positional embeddings



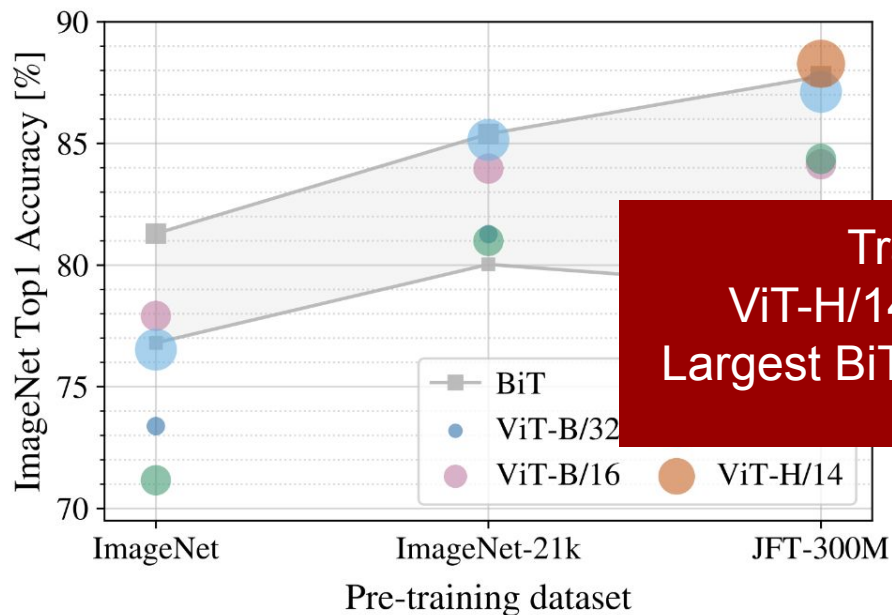
Vision Transformer



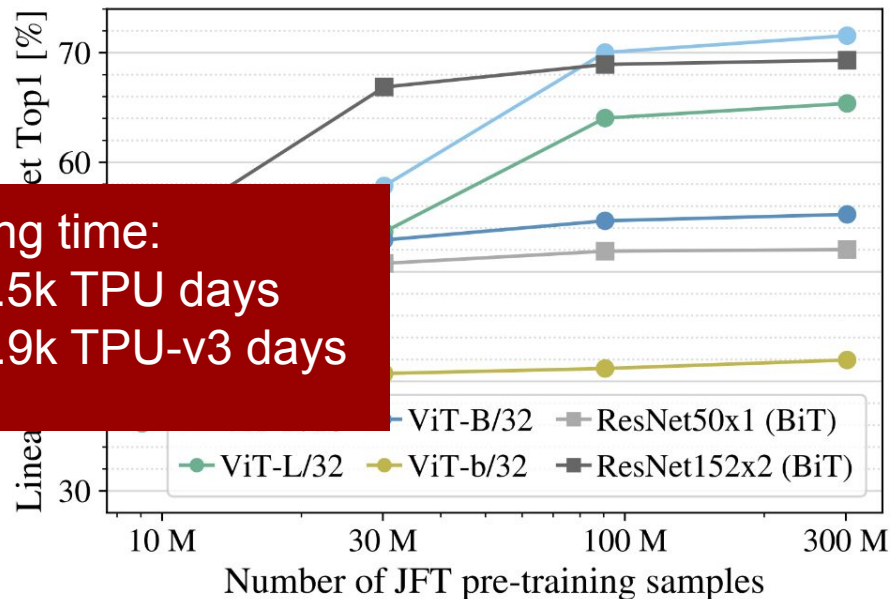
Transformer Encoder



ViT Results



Training time:
ViT-H/14: 2.5k TPU days
Largest BiT: 9.9k TPU-v3 days



ViT Summary

Model:

- Model is almost identical to Transformers for text sequences
- Replace words with $P \times P$ pixel image patches, $P \in \{14, 16, 32\}$ (no overlap)
- Each patch is embedded linearly into a vector of size 1024
- 1D **learned** positional embeddings

Training:

- For pre-training, optimize for image classification on large supervised dataset (e.g. ImageNet 21K, JFT -300M)
- For fine-tuning, learn a new classification head on a small dataset (e.g. CIFAR-100)

ACTIVITY: When do ViTs outperform CNNs, and vice versa?

Think about what you know about transformers -
what are some of their drawbacks?

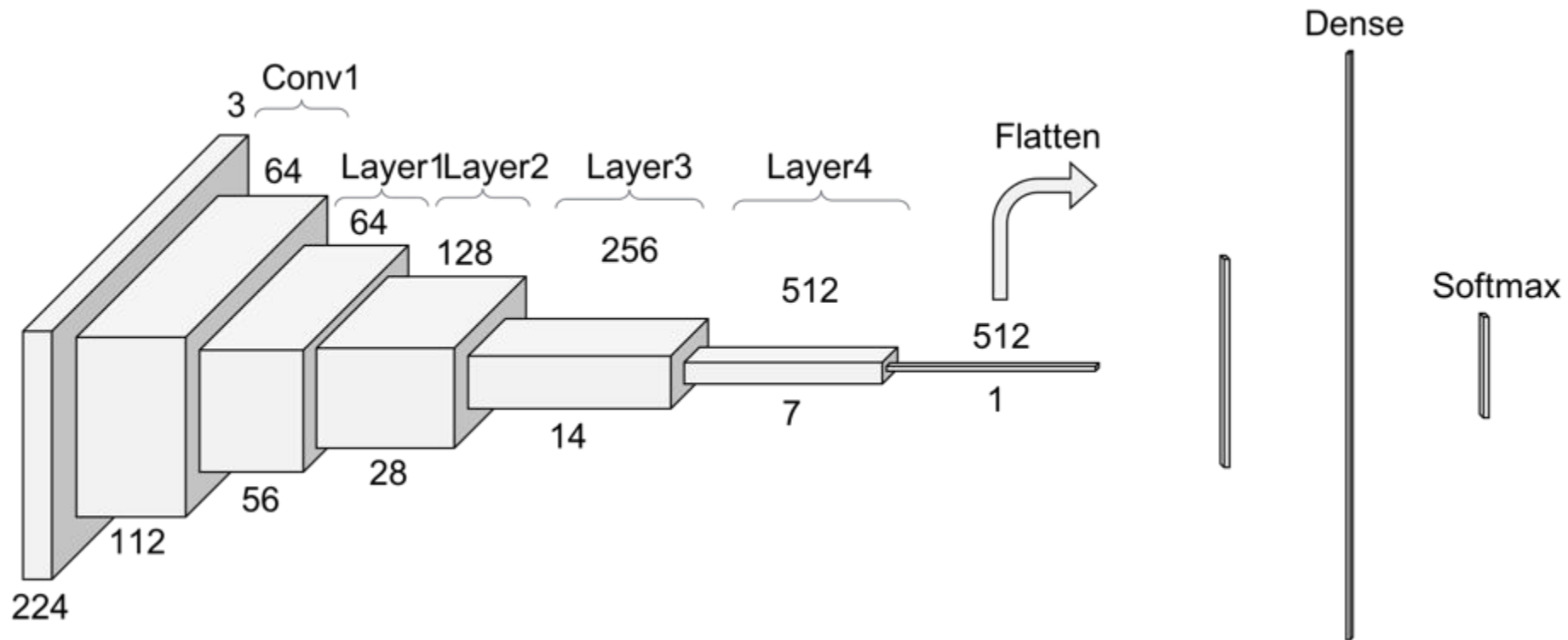
When is it “worth it” to use transformers instead of
just CNNs?

CNNs

- - Translational invariance
 - Simple, proven architecture
 - Capture features at different scales
 - Required less data

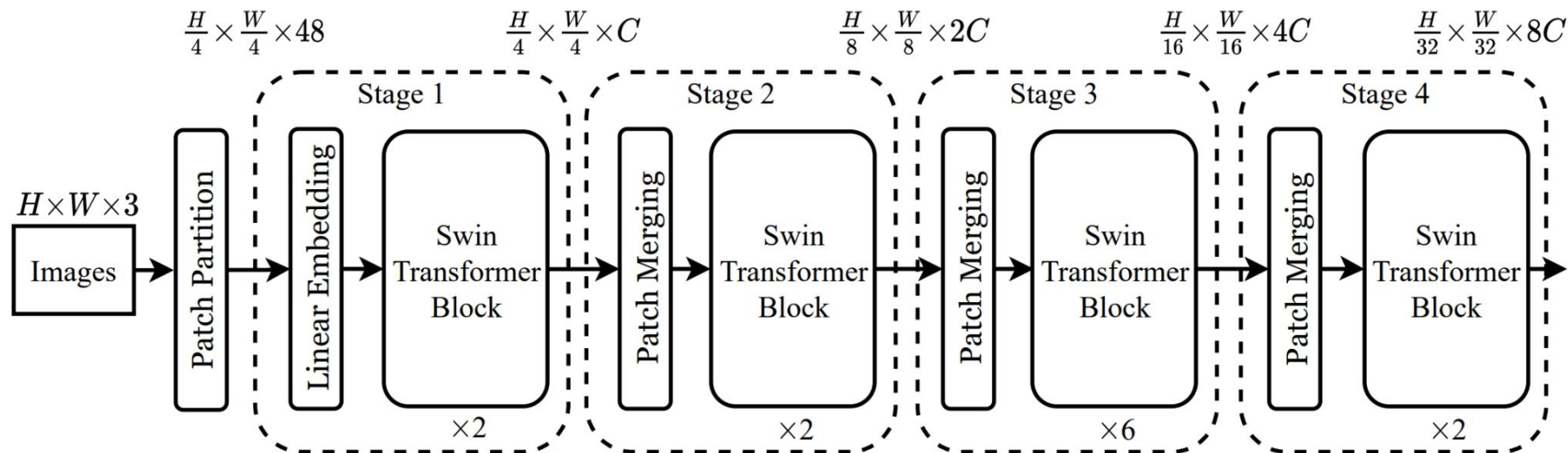
ViTs

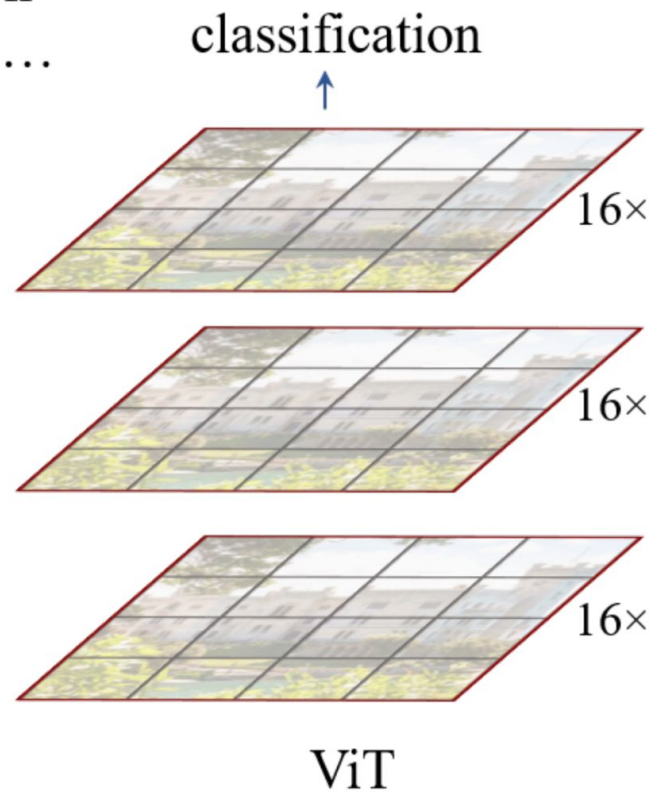
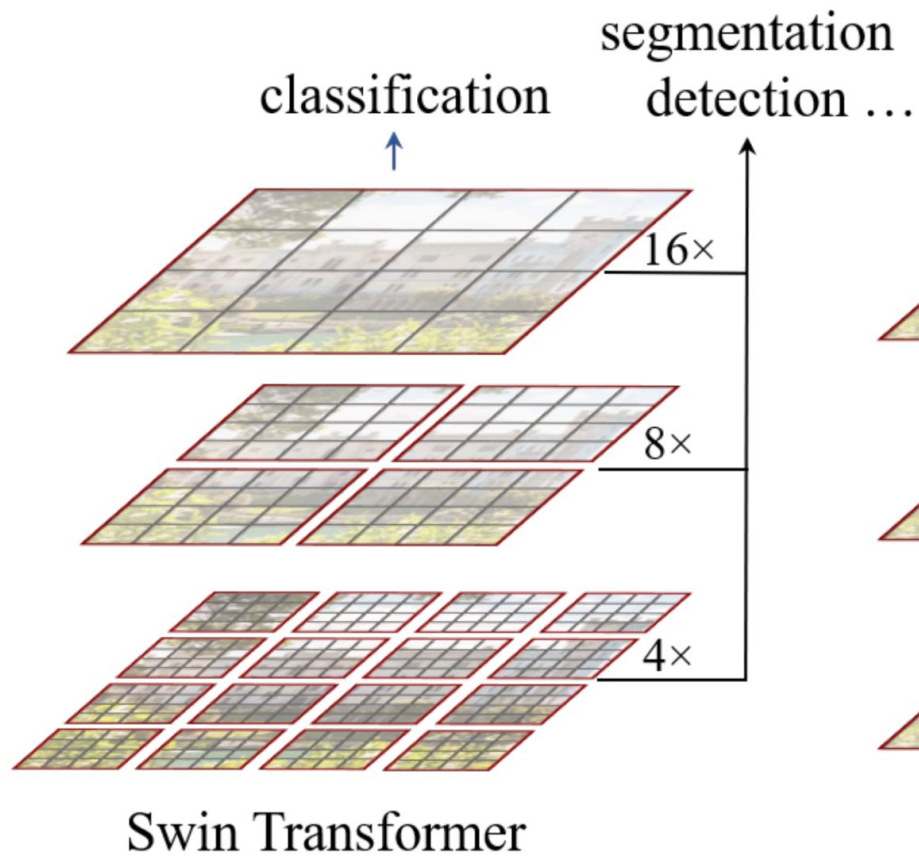
- Attention mechanisms
- Allow for multimodal data integration



Idea 3: Swim Transformers

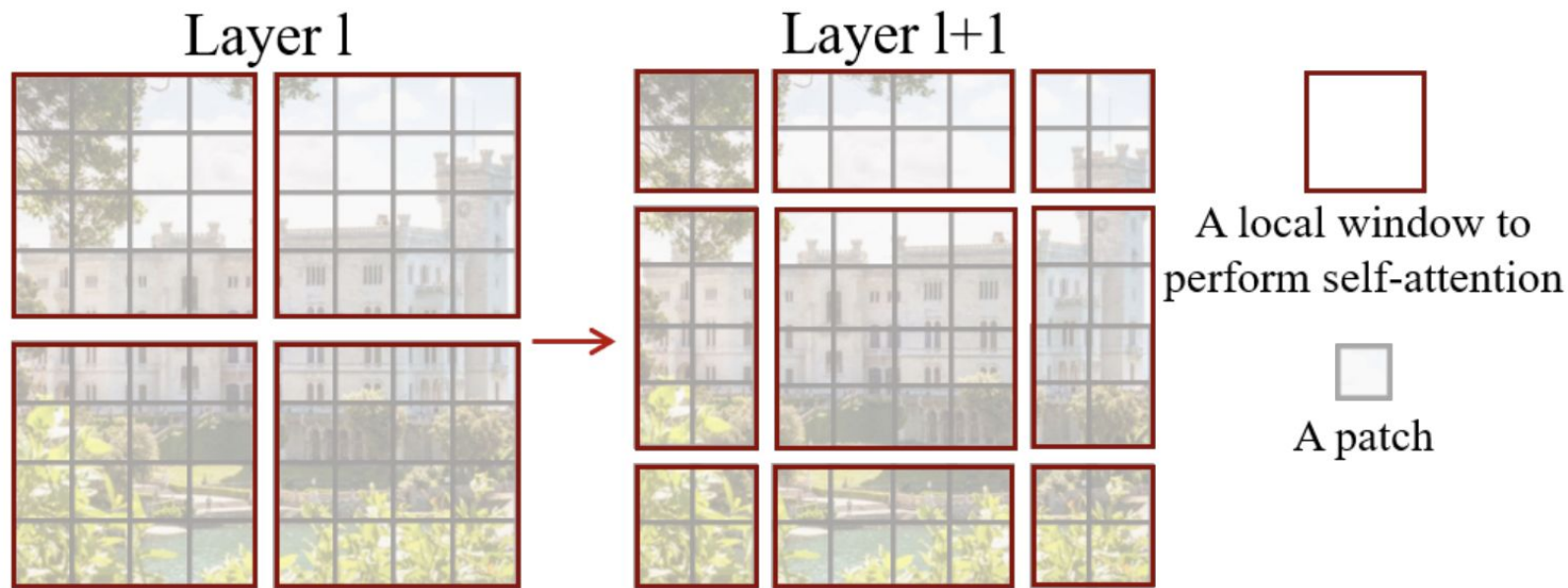
Hierarchical architecture that has the flexibility to model at various scales





Shifted Window attention

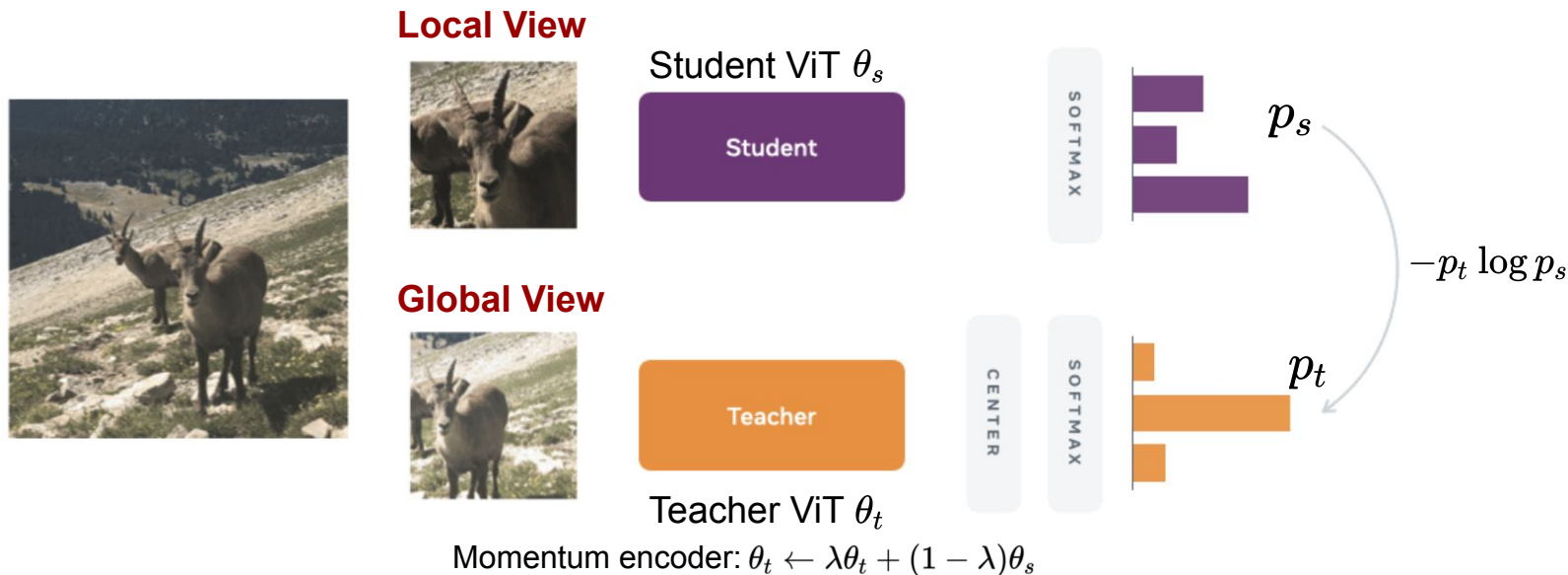
Linear computational complexity with respect to image size



Performance

method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
R-101x3 [38]	384^2	388M	204.6G	-	84.4
R-152x4 [38]	480^2	937M	840.5G	-	85.4
ViT-B/16 [20]	384^2	86M	55.4G	85.9	84.0
ViT-L/16 [20]	384^2	307M	190.7G	27.3	85.2
Swin-B	224^2	88M	15.4G	278.1	85.2
Swin-B	384^2	88M	47.0G	84.7	86.4
Swin-L	384^2	197M	103.9G	42.1	87.3

Self-Supervised Vision Transformers (DiNO)



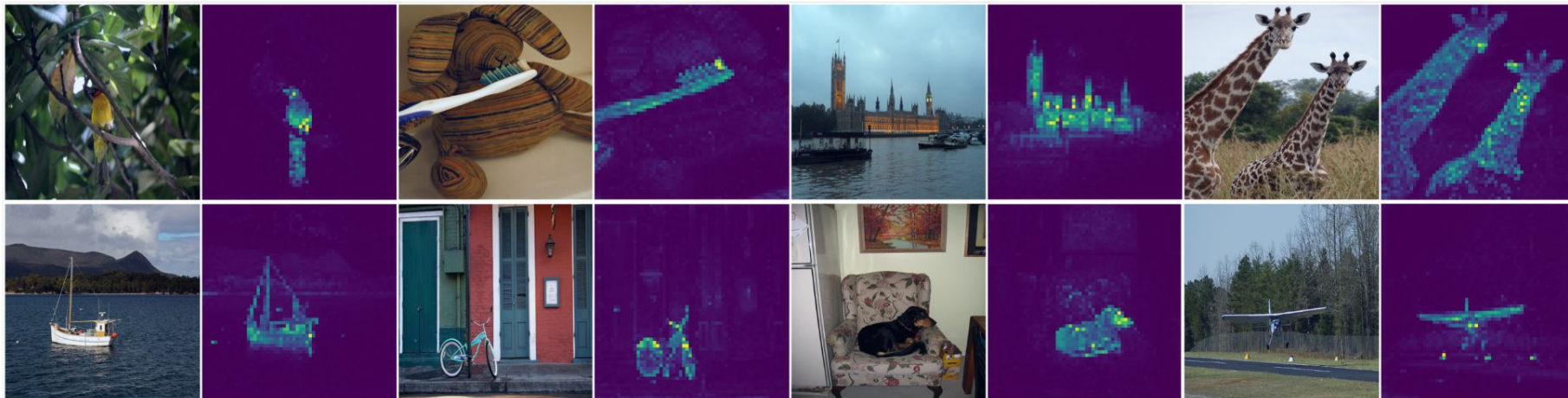
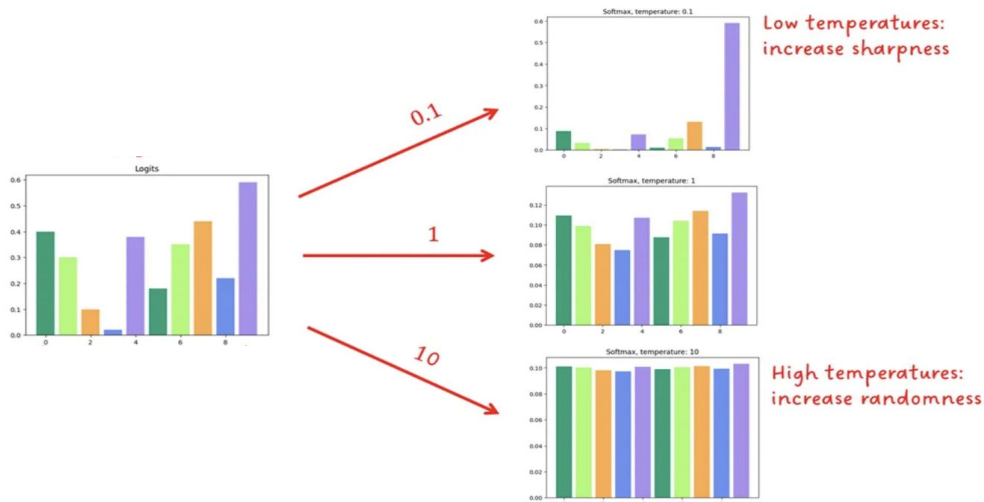


Figure 1: **Self-attention from a Vision Transformer with 8×8 patches trained with no supervision.** We look at the self-attention of the [CLS] token on the heads of the last layer. This token is not attached to any label nor supervision. These maps show that the model automatically learns class-specific features leading to unsupervised object segmentations.

Centering and sharpening

- Centering prevents one dimension from dominating
- Sharpening prevents learning a uniform distribution

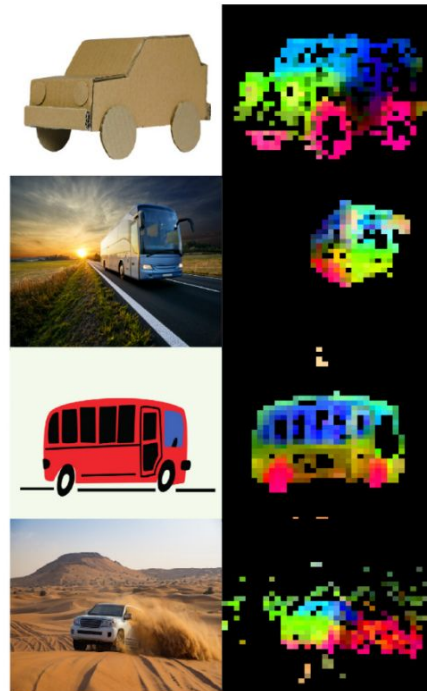
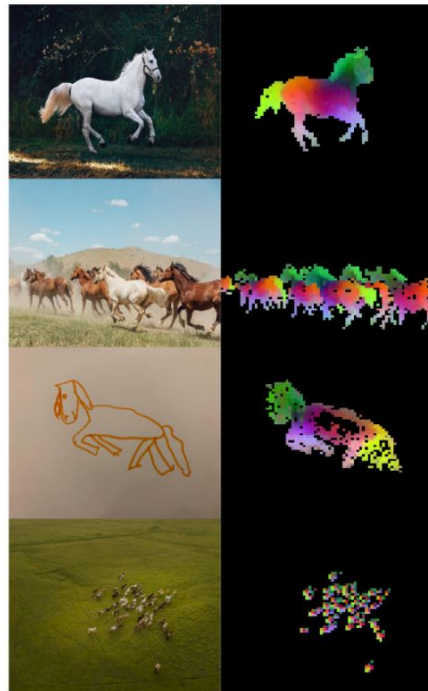
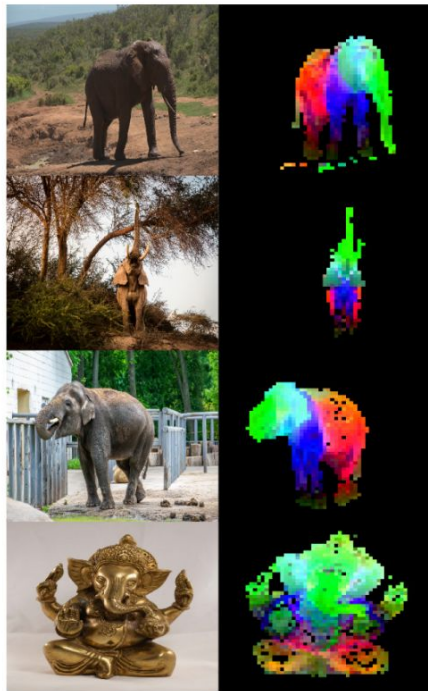
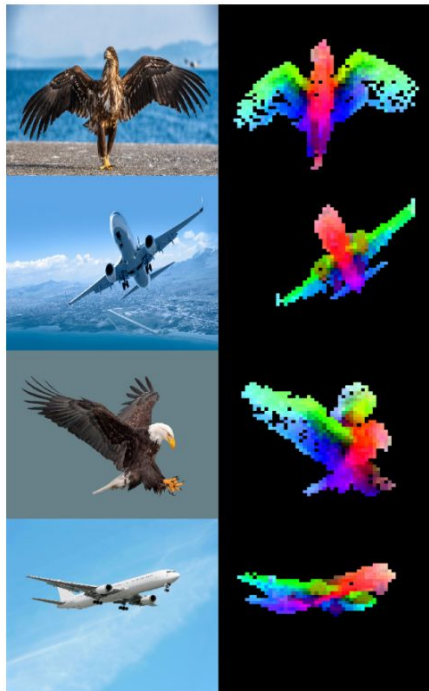
$$\text{softmax}(x)_i = \frac{e^{\frac{y_i}{T}}}{\sum_j^N e^{\frac{y_j}{T}}}$$



DINO v2

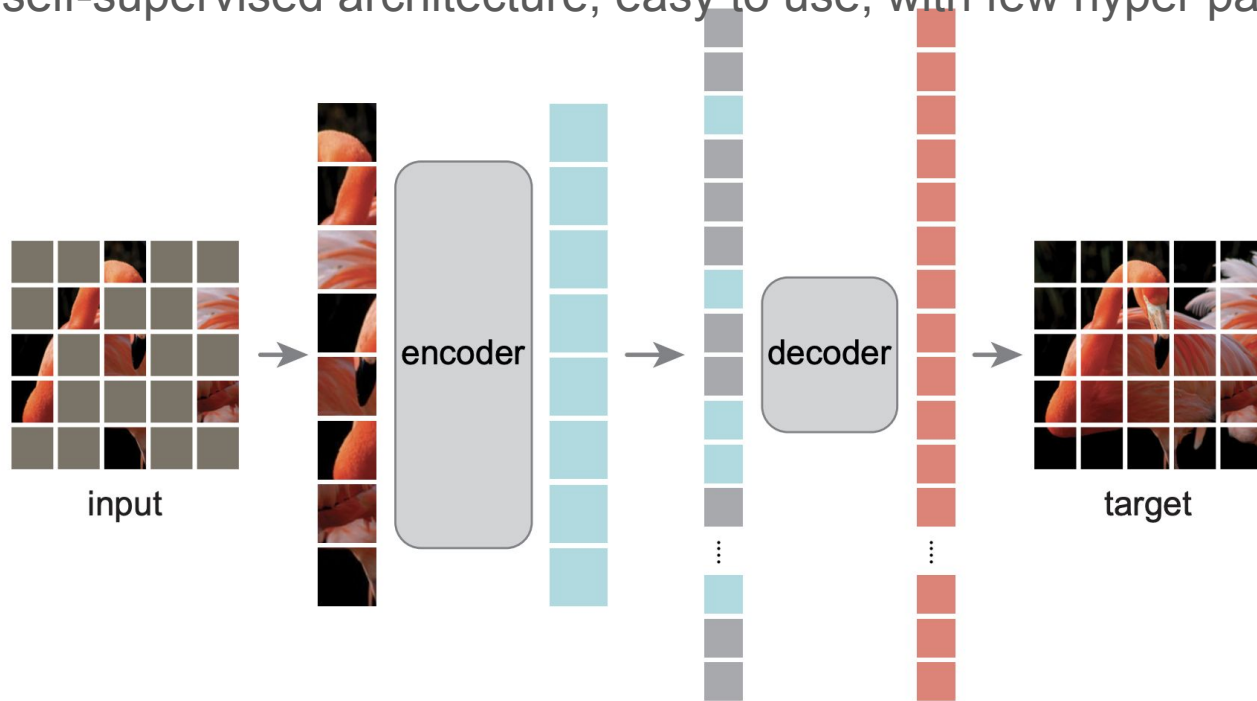
	INet-1k k-NN	INet-1k linear
iBOT	72.9	82.3
+ (our reproduction)	74.5 $\uparrow 1.6$	83.2 $\uparrow 0.9$
+ LayerScale, Stochastic Depth	75.4 $\uparrow 0.9$	82.0 $\downarrow 1.2$
+ 128k prototypes	76.6 $\uparrow 1.2$	81.9 $\downarrow 0.1$
+ KoLeo	78.9 $\uparrow 2.3$	82.5 $\uparrow 0.6$
+ SwiGLU FFN	78.7 $\downarrow 0.2$	83.1 $\uparrow 0.6$
+ Patch size 14	78.9 $\uparrow 0.2$	83.5 $\uparrow 0.4$
+ Teacher momentum 0.994	79.4 $\uparrow 0.5$	83.6 $\uparrow 0.1$
+ Tweak warmup schedules	80.5 $\uparrow 1.1$	83.8 $\uparrow 0.2$
+ Batch size 3k	81.7 $\uparrow 1.2$	84.7 $\uparrow 0.9$
+ Sinkhorn-Knopp	81.7 =	84.7 =
+ Untying heads = DINOv2	82.0 $\uparrow 0.3$	84.5 $\downarrow 0.2$

Cornell Bowers C-IS

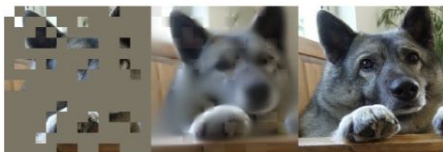
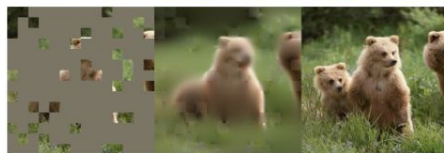
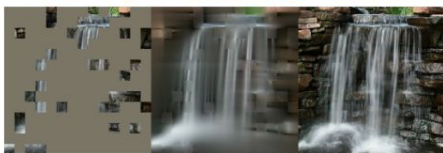
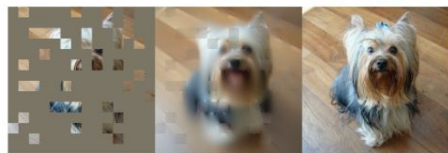
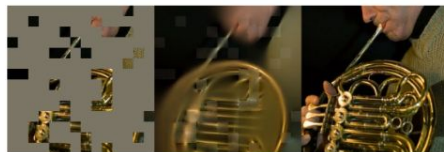
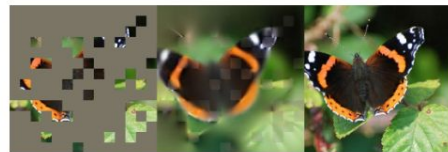
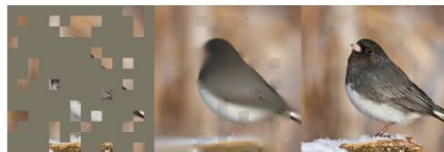
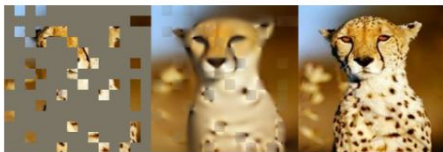


Masked Autoencoders (MaE)

A simple self-supervised architecture, easy to use, with few hyper parameters.



Cornell Bowers C-IS



Discuss: BERT is trained with cross entropy loss. What loss function should we use for MaE?

MaE Results

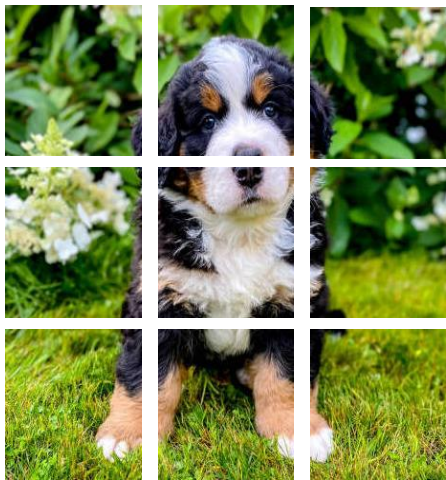
- Compared to supervised ViTs
 - Requires minimal data augmentation
 - Transfers better to downstream vision tasks
 - Object detection, segmentation

case	ft	lin
none	84.0	65.7
crop, fixed size	84.7	73.1
crop, rand size	84.9	73.5
crop + color jit	84.3	71.9

(e) **Data augmentation.** Our MAE works with minimal or no augmentation.

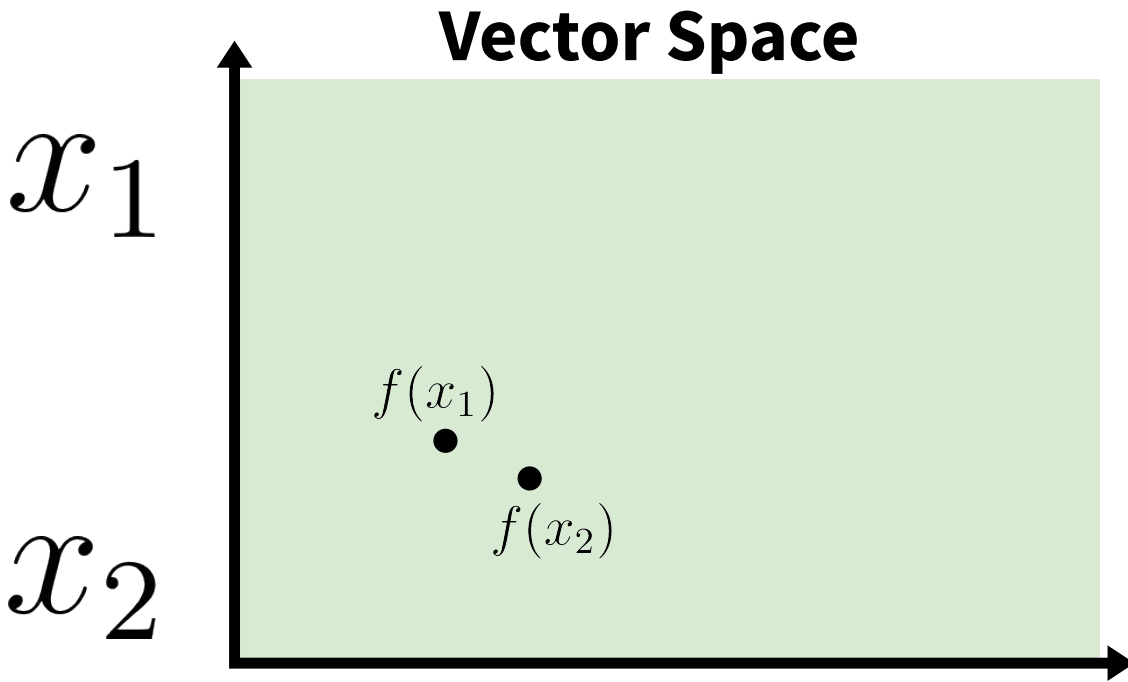
method	pre-train data	AP ^{box}		AP ^{mask}	
		ViT-B	ViT-L	ViT-B	ViT-L
supervised	IN1K w/ labels	47.9	49.3	42.9	43.9
MoCo v3	IN1K	47.9	49.3	42.7	44.0
BEiT	IN1K+DALLE	49.8	53.3	44.4	47.1
MAE	IN1K	50.3	53.3	44.9	47.2

Table 4. **COCO object detection and segmentation** using a ViT Mask R-CNN baseline. All entries are based on our implementation. Self-supervised entries use IN1K data *without* labels. Mask AP follows a similar trend as box AP.

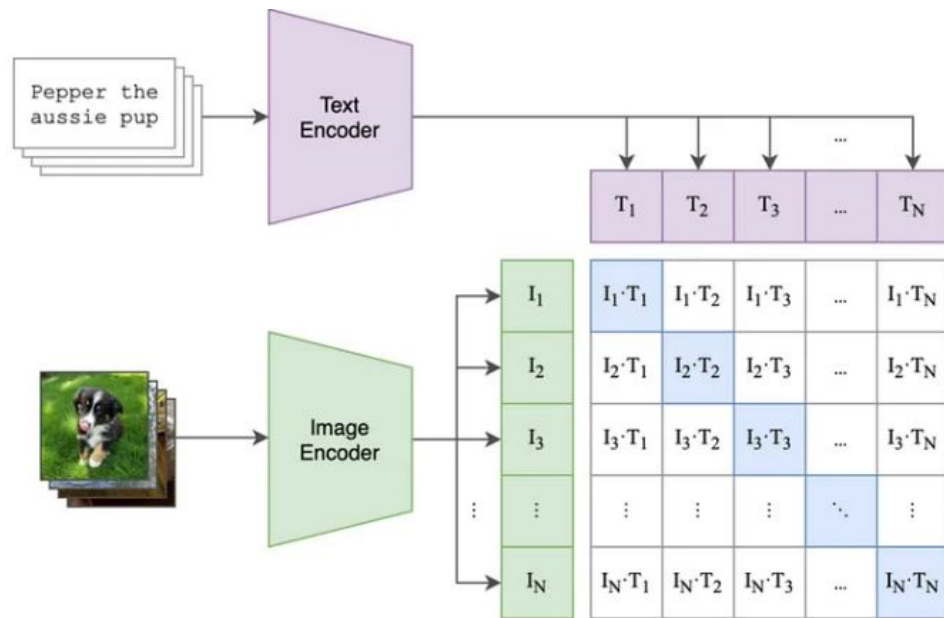


A puppy sits
with flowers

$f(x)$ = transformer rep.



CLIP (Contrastive Language–Image Pre-training)



Discuss: How can you train this model?

```
# image_encoder - ResNet or Vision Transformer
# text_encoder  - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, l]       - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t             - learned temperature parameter

# extract feature representations of each modality
I_e = image_encoder(I) # [n, d_i]
```

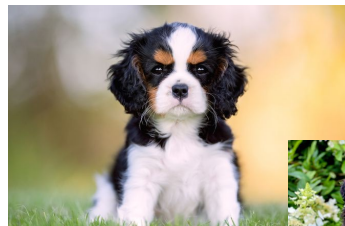
- Trained on 256 V100 GPUs for two weeks on 400 million (image, text pairs)
- On AWS, this would cost at least 200k dollars

```
# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

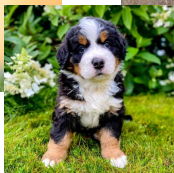
# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss   = (loss_i + loss_t)/2
```

Ranking using CLIP

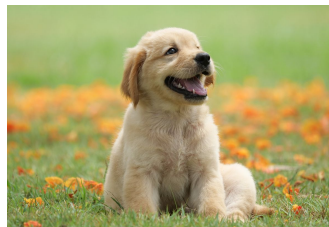
sim score - 0.9



sim score - 0.3



sim score - 0.6



sim score - 0.1



Vector Space

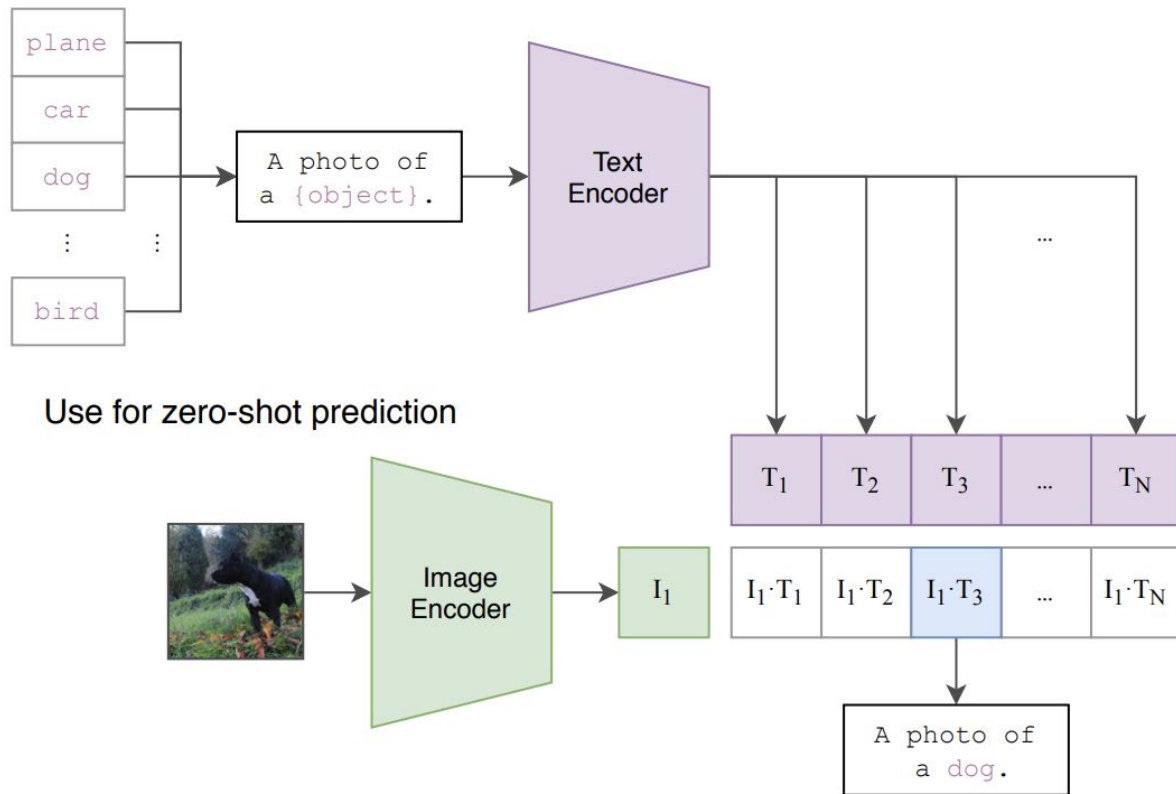


Cornell Bowers C-IS

Clip demo

<https://huggingface.co/spaces/vivien/clip>

Create dataset classifier from label text



Cornell Bowers C-IS

FOOD101

guacamole (90.1%) Ranked 1 out of 101 labels



✓ 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.

YOUTUBE-BB

airplane, person (89.0%) Ranked 1 out of 23



✓ a photo of **airplane**.

✗ a photo of **bird**.

✗ a photo of **bear**.

✗ a photo of **giraffe**.

✗ a photo of **car**.

SUN397

television studio (90.2%) Ranked 1 out of 397



✓ a photo of a **television studio**.

✗ a photo of a **podium indoor**.

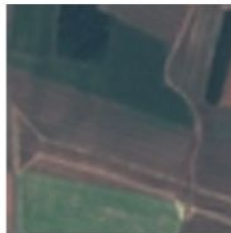
✗ a photo of a **conference room**.

✗ a photo of a **lecture room**.

✗ a photo of a **control room**.

EUROSAT

annual crop land (12.9%) Ranked 4 out of 10



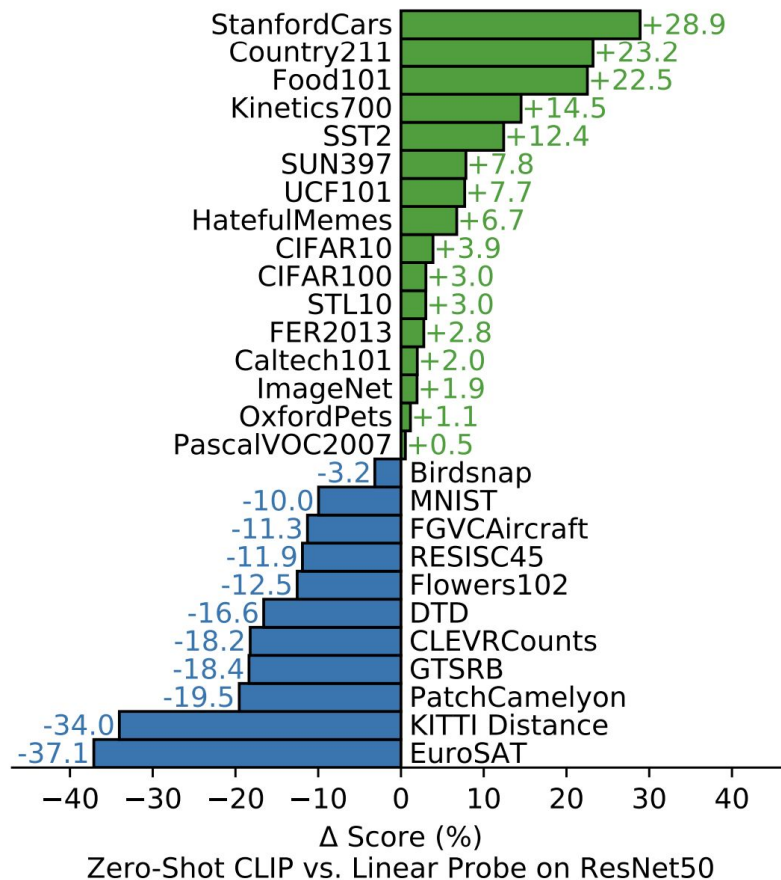
✗ a centered satellite photo of **permanent crop land**.

✗ a centered satellite photo of **pasture land**.

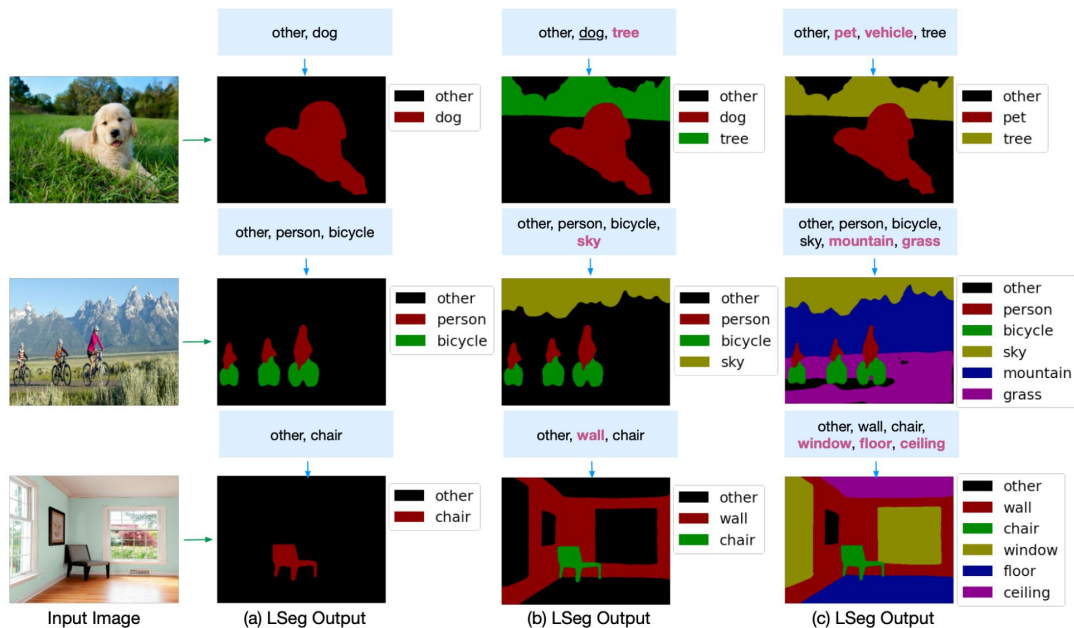
✗ a centered satellite photo of **highway or road**.

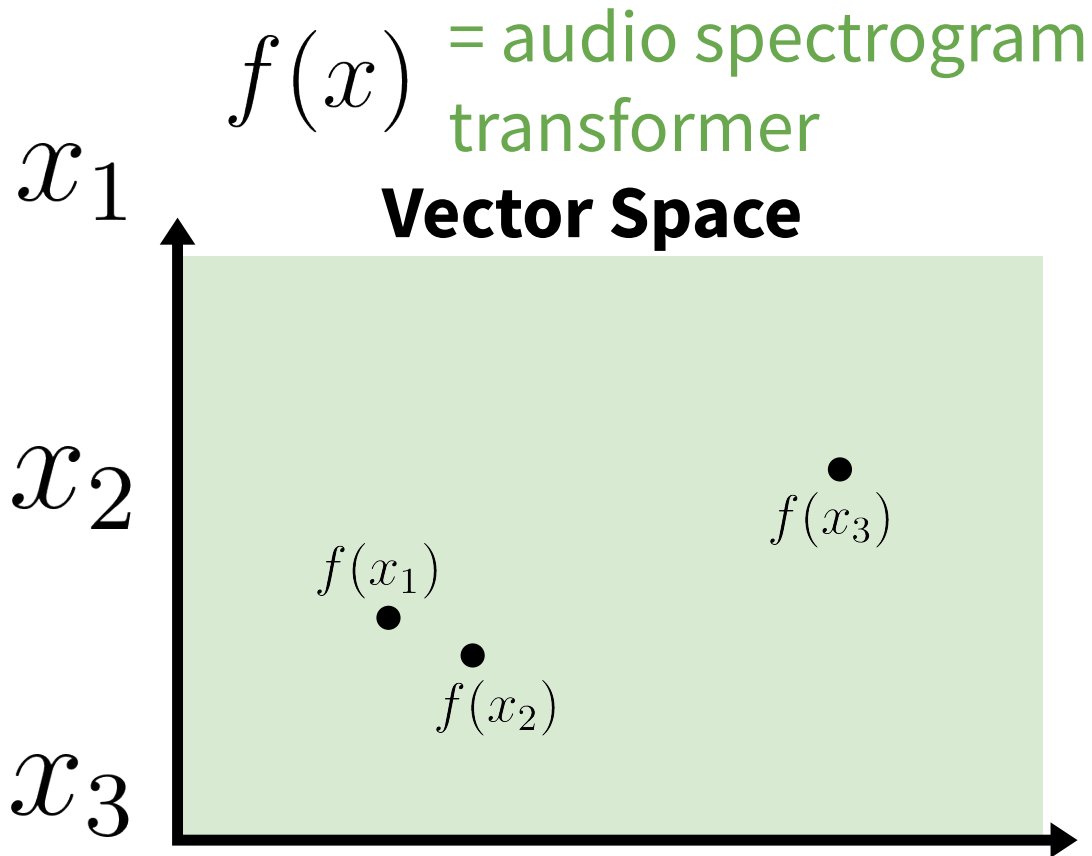
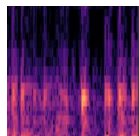
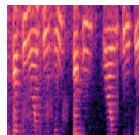
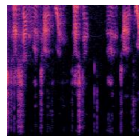
✓ a centered satellite photo of **annual crop land**.

✗ a centered satellite photo of **brushland or shrubland**.

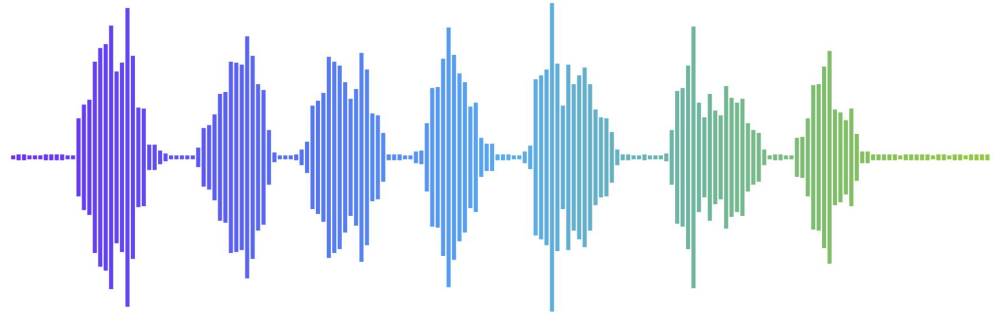


Application of CLIP





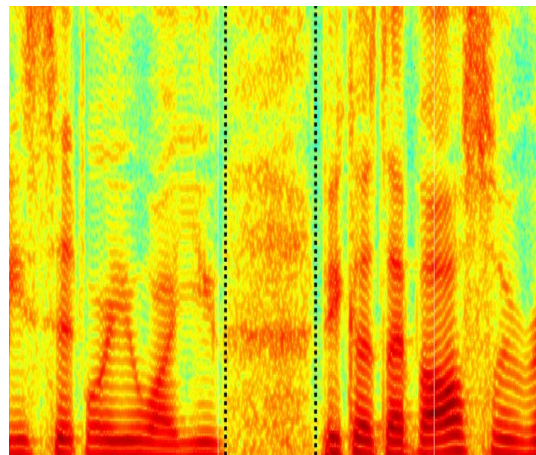
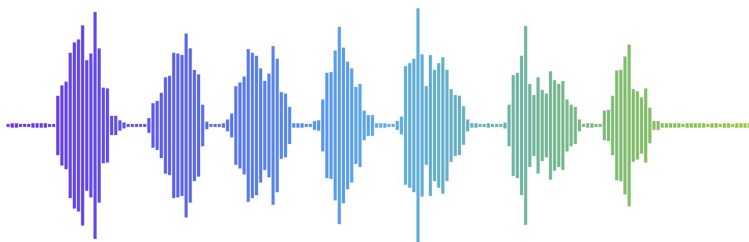
Audio Processing



Audio File

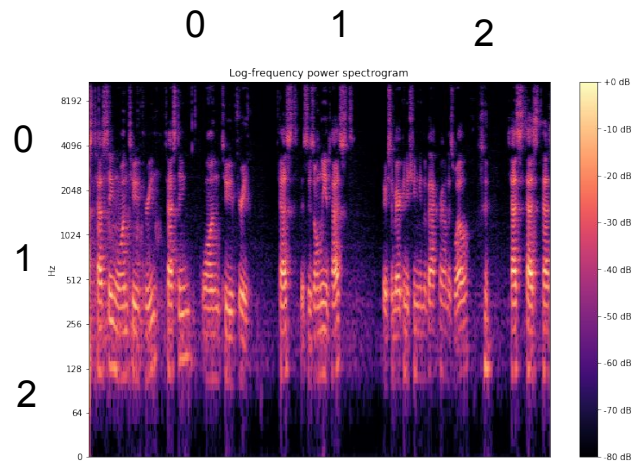
Spectrogram:

- Energy, pitch, fundamental frequency
- Decomposes signal into frequencies and their corresponding amplitudes

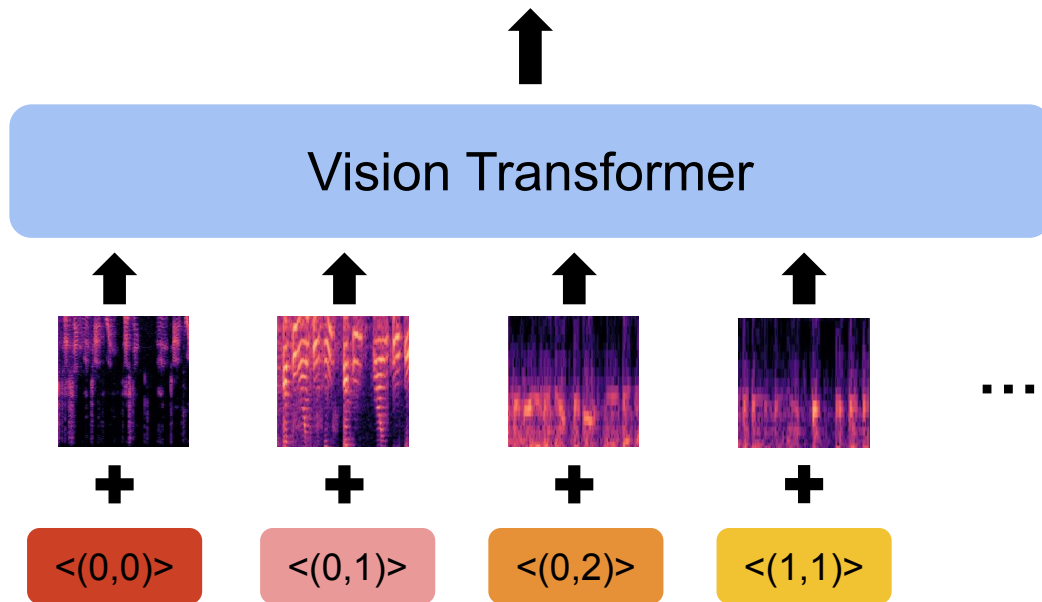


Spectrogram source: [Dumpala 2017](#)

Audio as a vision problem



Most likely word sequence



Review

- Transformers can be used for vision tasks
 - Typically do not use convolution or other image specific architectures - just flattened image patches
 - Consequently they require more data
 - However they are very flexible and can be used with other data modalities
 - Learn positional embedding
- Swin transformers
 - Learn features at different scales
 - Scale linearly with the image size
 - Are less popular because of their complexity
- Self-supervised learning
 - Provides Transformer based models with the amount of data they need
 - Often used as backbone pre-trained models
 - Dino and MAE both learn very good embeddings
- Multi-modality
 - CLIP learns a common embedding for images and text
 - Can be used to retrieve images or parts of images from plain text