

Cornell Bowers CIS

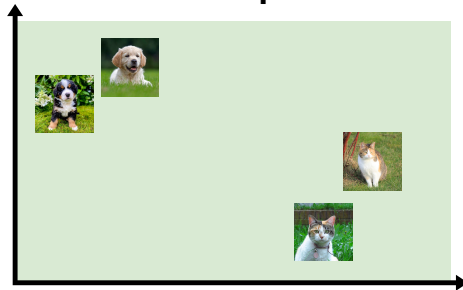
Logistics

- HW3 is due next week
- We have a feedback form (due Friday, March 8th)
- Project proposal due today
 - You can use slip days
 - Please think about compute when you choose a paper for your project
 - Look for reported training hardware and training times
 - Choose papers from ICML, ICLR, NeurIPS, CVPR, ECCV, ACL, EMNLP, NAACL

Cornell Bowers CIS

Use self-supervised learning to learn embeddings for images

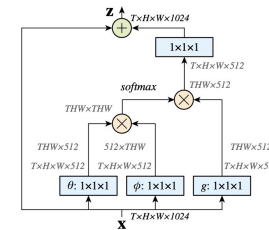
Vector Space



Cornell Bowers CIS

How to use Attention for Vision Tasks?

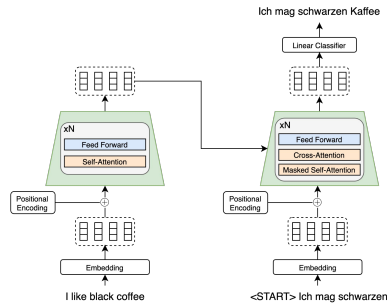
Idea #1: Add attention to existing CNNs



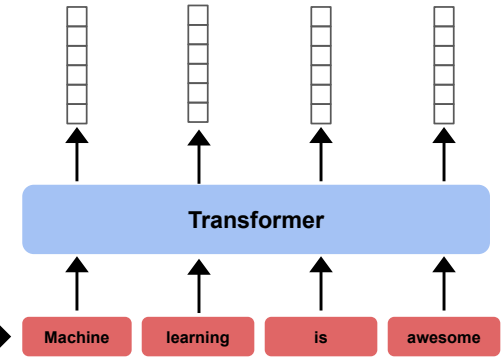
Wang, X., Girshick, R., Gupta, A., & He, K. (2018). Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 7794-7803).

How to use Attention for Vision Tasks?

Idea #2: Adapt standard transformers to image data

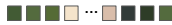
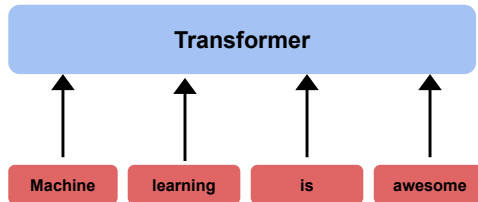
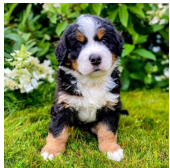


Machine learning
is awesome

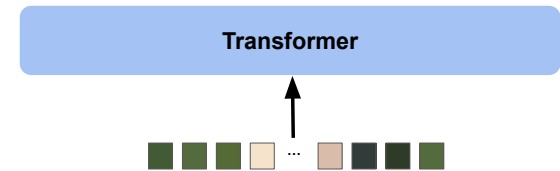
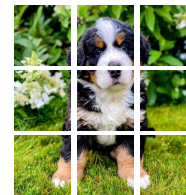


Can we extend this idea to images?

Machine learning
is awesome

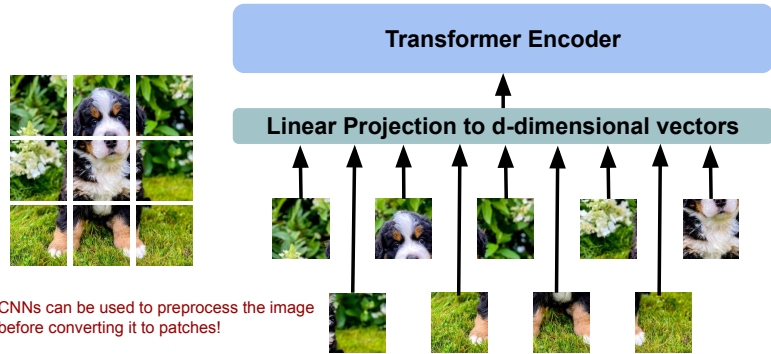


Use pixels as input

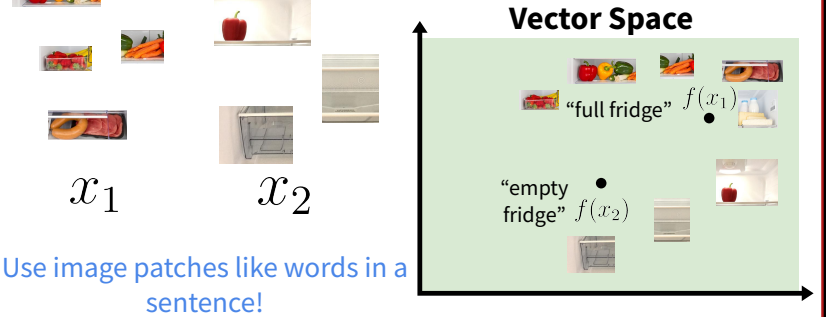


$R \times R$ image needs R^2 elements per attention matrix
 $R=128$, 48 layers, 16 heads per layer takes 768GB of memory for attention matrices for a single example...

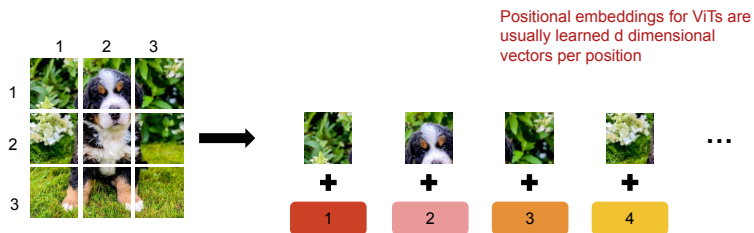
Image patches as “words”



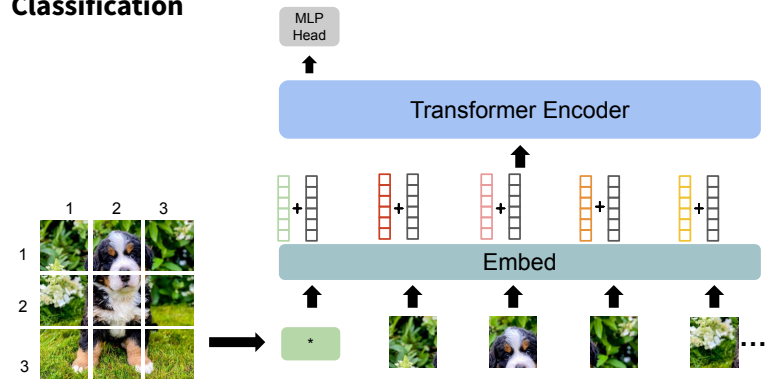
$$f(x) = \text{Vision Transformers}$$



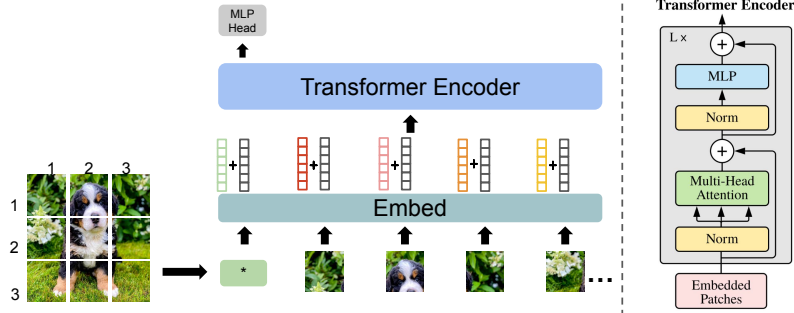
Adding positional embeddings



Classification

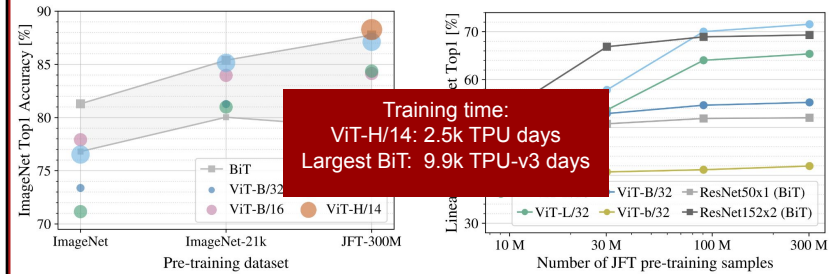


Vision Transformer



Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.

ViT Results



Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.

ViT Summary

Model:

- Model is almost identical to BERT
- Replace words with PxP pixel image patches, $P \in \{14, 16, 32\}$ (no overlap)
- Each patch is embedded linearly into a vector of size 1024
- 1D 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)

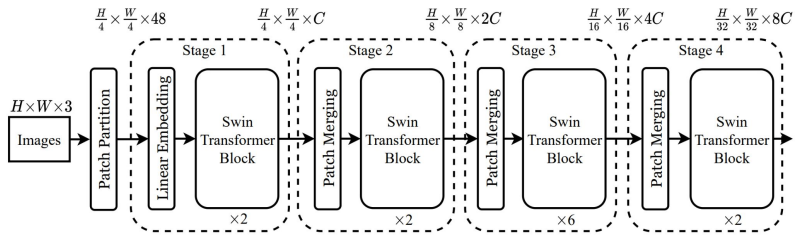
ACTIVITY: When do ViTs outperform CNNs, and vice versa?

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

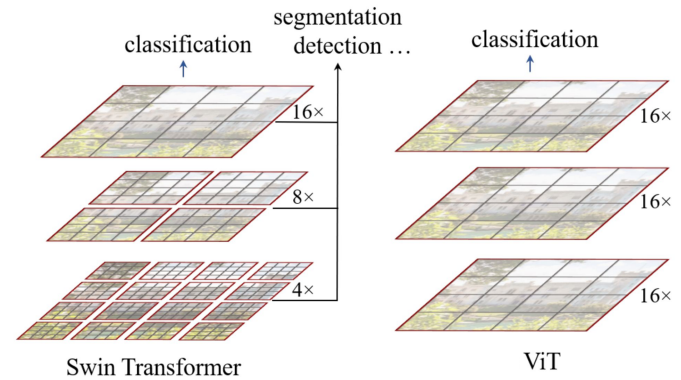
When is it “worth it” to use transformers instead of just CNNs?

Swim Transformers

Hierarchical architecture that has the flexibility to model at various scales

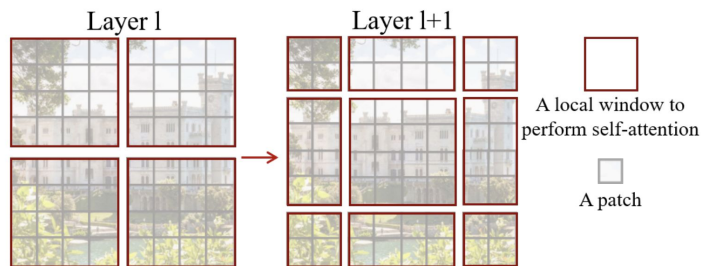


Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., ... & Guo, B. (2021). Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision* (pp. 10012-10022).



Shifted Window attention

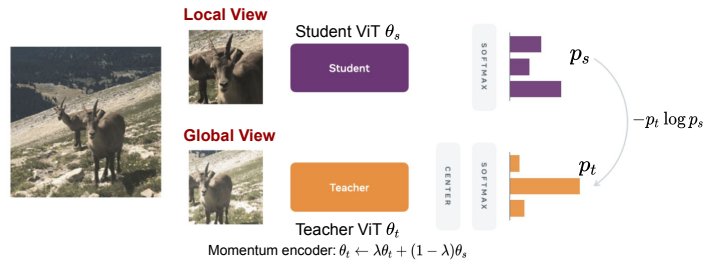
Linear computational complexity with respect to image size



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)



Centering and sharpening

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

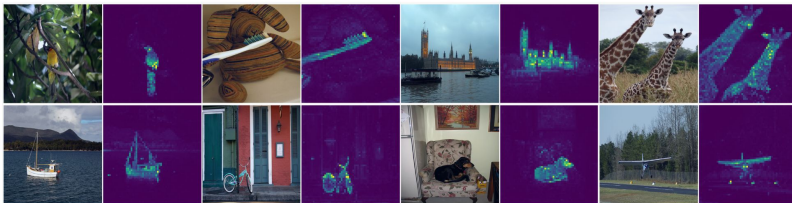
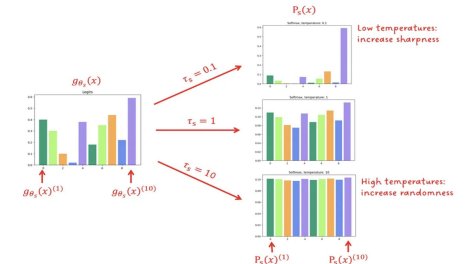
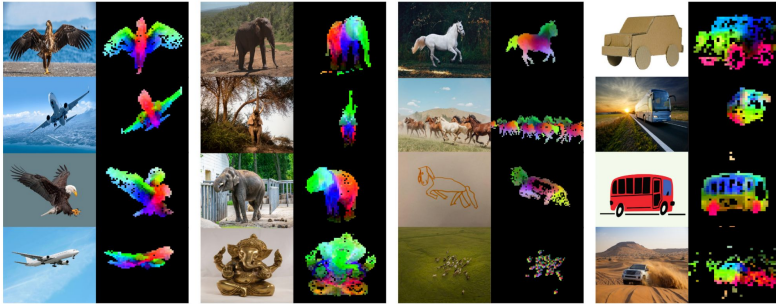


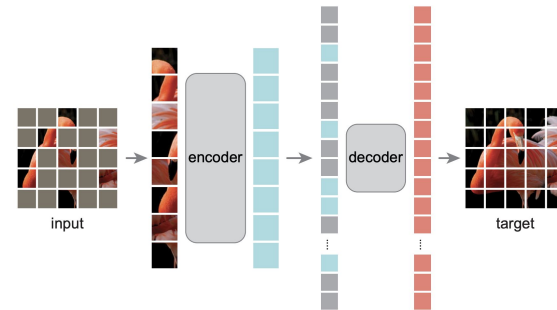
Figure 1: **Self-attention from a Vision Transformer with 8×8 patches trained with no supervision.** We look at the self-attention of the [CLS] token on the heads of the last layer. This token is not attached to any label nor supervision. These maps show that the model automatically learns class-specific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

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

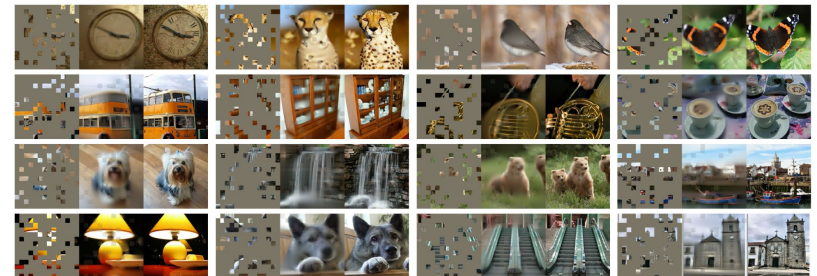


Masked Autoencoders (MaE)



He, K., Chen, X., Xie, S., Li, Y., Dollár, P., & Girshick, R. (2022). Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 16000-16009).

Discuss: BERT is trained with cross entropy loss. Can you do the same with MaE or should you use a different loss?



MaE Results

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

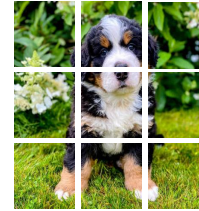
case	ft	lin
none	84.0	65.7
crop, fixed size	84.7	73.1
crop, rand size	84.9	73.5
crop + color jit	84.3	71.9

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

method	pre-train data	A ^{pbox}		A ^{pmask}	
		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
BET	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.

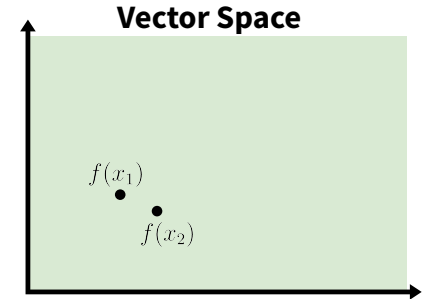
$$f(x) = \text{transformer rep.}$$



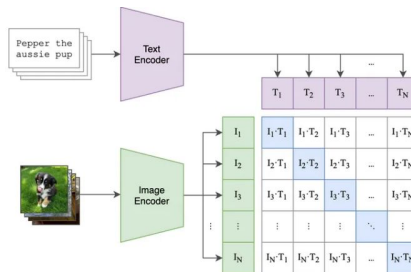
A puppy sits with flowers

x_1

x_2



CLIP (Contrastive Language-Image Pre-training)



Conde, M. V., & Turgutlu, K. (2021). CLIP-Art: Contrastive pre-training for fine-grained art classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 3956-3960).

Discuss: How can you train this model?

Cornell Bowers CIS

```
# 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_f = image_encoder(I) #[n, d_i]

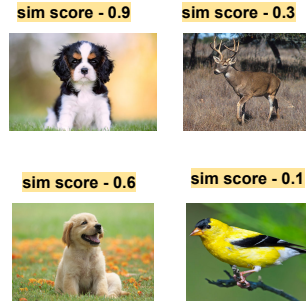
# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

# 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
```

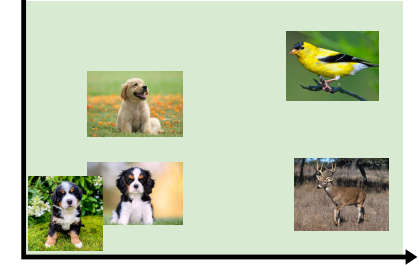
- Trained on 256 V100 GPUs for two weeks on 400 million (image, text pairs)
- On AWS, this would cost at least 200k dollars

Cornell Bowers CIS

Ranking using CLIP



Vector Space



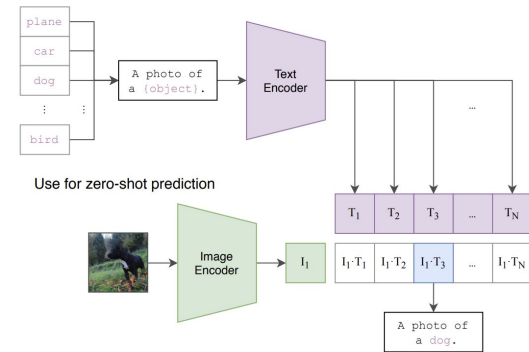
Cornell Bowers CIS

Clip demo

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

Cornell Bowers CIS

Create dataset classifier from label text



Cornell Bowers CIS

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 a **airplane**.
- ✗ a photo of a bird.
- ✗ a photo of a bear.
- ✗ a photo of a giraffe.
- ✗ a photo of a car.

SUN397

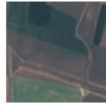
television studio (90.2%) Ranked 1 out of 397



- ✓ a photo of a **television studio**.
- ✗ a photo of a podium indoors.
- ✗ a photo of a conference room.
- ✗ a photo of a lecture room.
- ✗ a photo of a control room.

EUROSAT

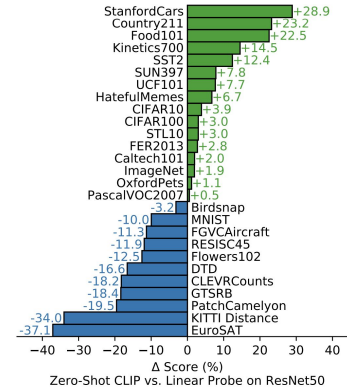
annual crop land (12.9%) Ranked 4 out of 10



- ✗ 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.

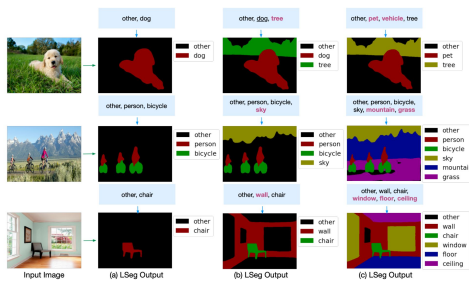
<https://openai.com/research/clip>

Cornell Bowers CIS



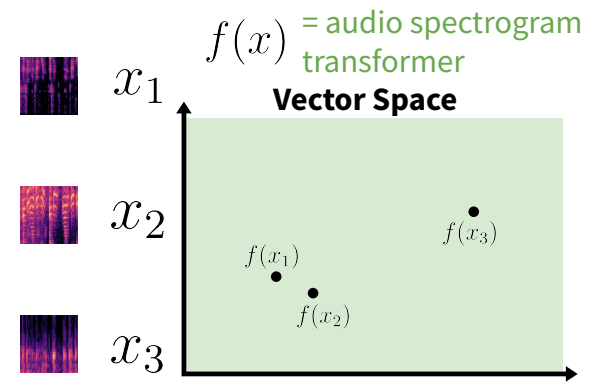
Cornell Bowers CIS

Application of CLIP



<https://arxiv.org/pdf/2201.03546.pdf>

Cornell Bowers CIS

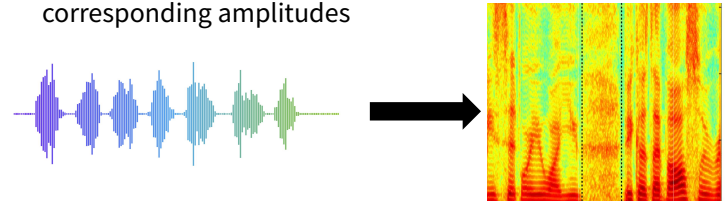


Audio Processing



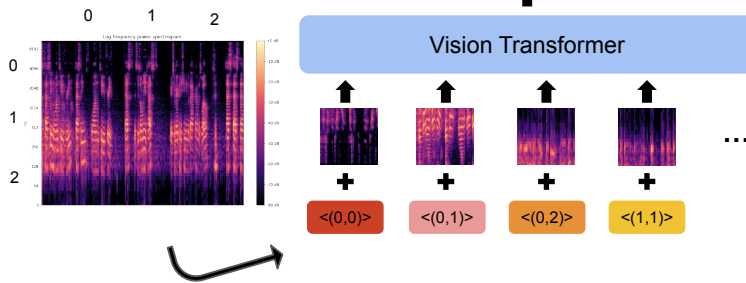
Spectrogram:

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



Audio as a vision problem

Most likely word sequence



Review

- Transformers can be used for vision tasks
- Swin transformers can be used for learning features at different scales
- Self-supervised learning is also helpful for transformer backbone vision models
 - Dino and MAE both learn very good embeddings
- Using transformer models for images and text helps build multi-modal models like CLIP