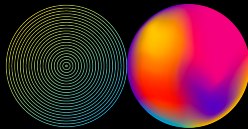# Recent Trends in Machine Learning:
# A Large-scale Perspective

## A Short Introduction to **Multi-modal AI** Models (Part 1)

**Saehoon Kim @ Kakaobrain**

# Outline of This Course

**CLIP**
**Encoder-only**

**DALL-E**
**Decoder-only**

**DALL-E 2**
**Enc-Dec**

**05/04**

**05/11**

**05/18**

# Outline of This Course

**Contrastive Learning**

**DALL-E**
Decoder-only

**DALL-E 2**
Enc-Dec

# Outline of This Course



presets | a girl running on the beach

**DALL-E**
**Decoder-only**

**DALL-E 2**
**Enc-Dec**

# Outline of This Course

**Contrastive Learning**

**Autoregressive Model**

**DALL-E 2**
**Enc-Dec**

# Outline of This Course

**Contrastive Learning**



Autoregressive Model

**DALL-E 2**
**Enc-Dec**

# Outline of This Course
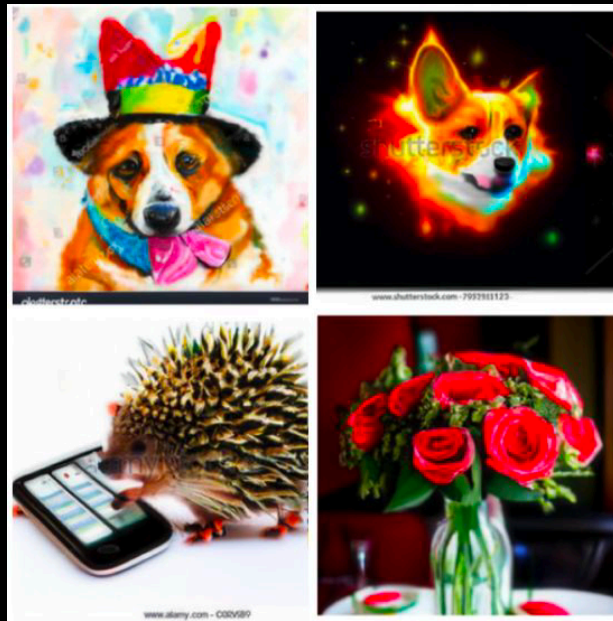
**Contrastive Learning**

**Autoregressive Model**

**Diffusion Model**

# Outline of This Course



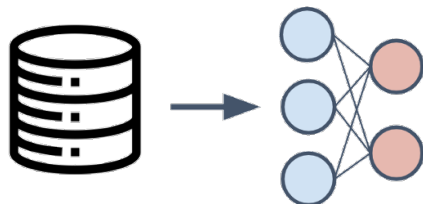**Contrastive Learning**

**Autoregressive Model**

# Background

**Self-Supervised Representation Learning**

# Transfer Learning
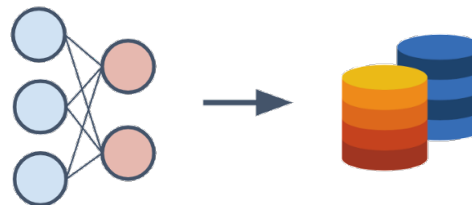
Transferring visual features learned from a large annotated set into small-scale downstream tasks has been significantly improved the performance!
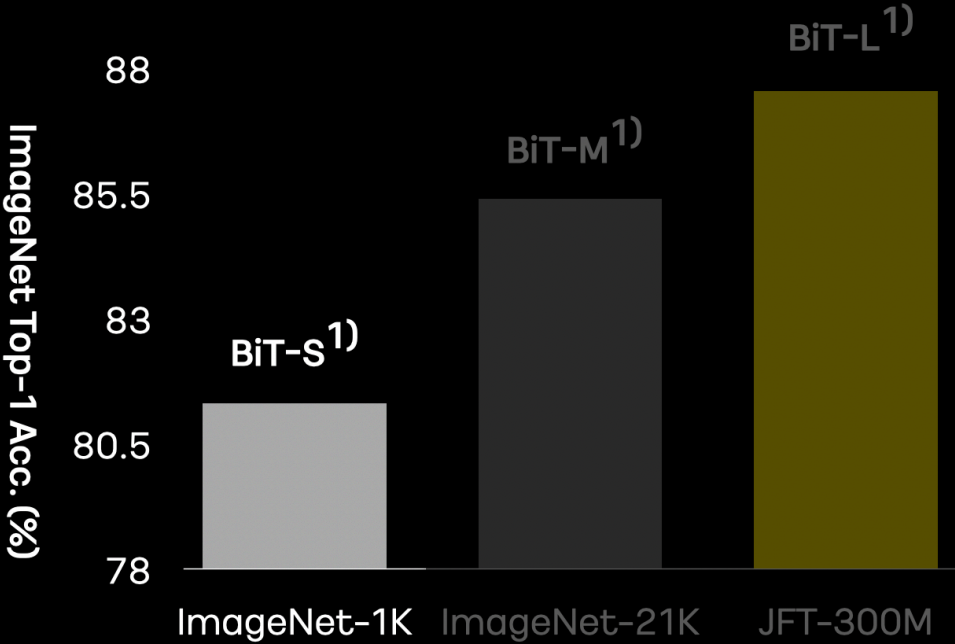
Kolesnikov et al., "Big Transfer (BiT): General Visual Representation Learning", ECCV'20. 10

# Transfer Learning



Kolesnikov et al., "Big Transfer (BiT): General Visual Representation Learning", ECCV'20.

# Transfer Learning



Kolesnikov et al., "Big Transfer (BiT): General Visual Representation Learning", ECCV'20.

# Transfer Learning



Kolesnikov et al., "Big Transfer (BiT): General Visual Representation Learning", ECCV'20.

# Transfer Learning



**Can we learn visual features without labeled samples in the upstream pre-training?**

Kolesnikov et al., "Big Transfer (BiT): General Visual Representation Learning", ECCV'20.

# Contrastive Learning

Learning the global representations by comparing the semantically similar and dissimilar images without human annotations
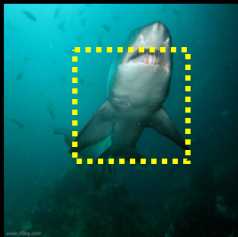
# Contrastive Learning

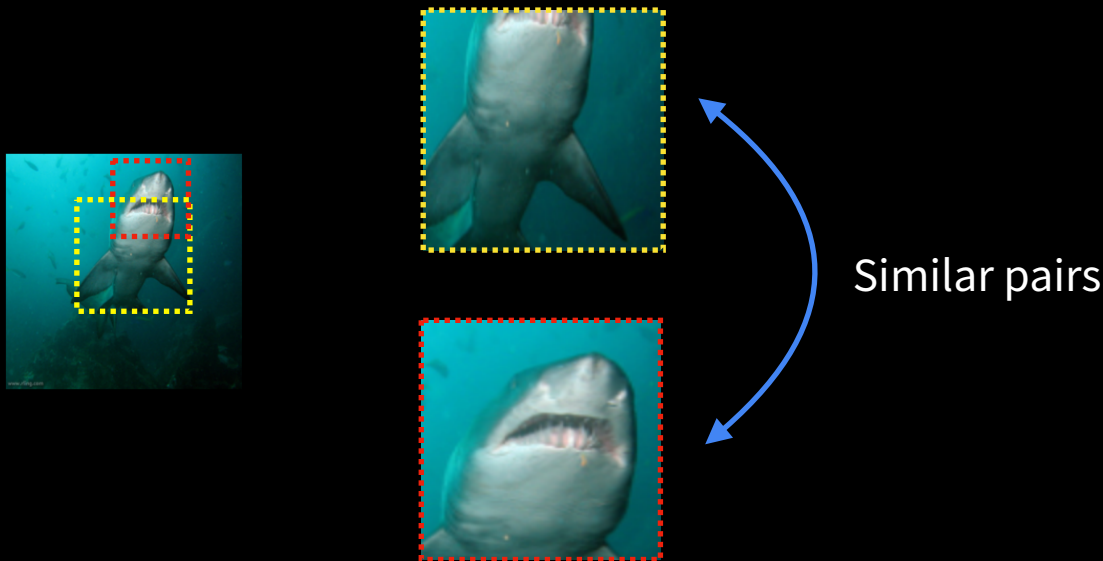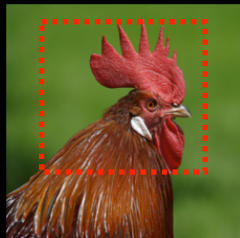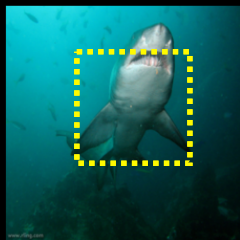How to automatically obtain similar and dissimilar pairs without labels?

# Contrastive Learning

How to automatically obtain similar and dissimilar pairs without labels?
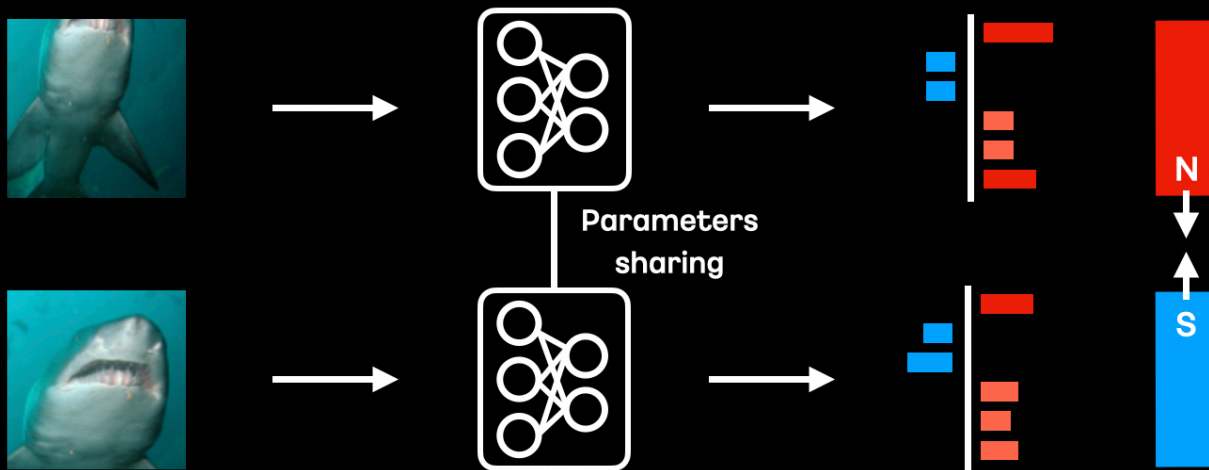
# Contrastive Learning

How to automatically obtain similar and dissimilar pairs without labels?



Similar pairs

# Contrastive Learning

How to automatically obtain similar and dissimilar pairs without labels?



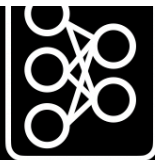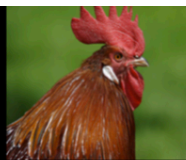Dissimilar pairs

# Contrastive Learning

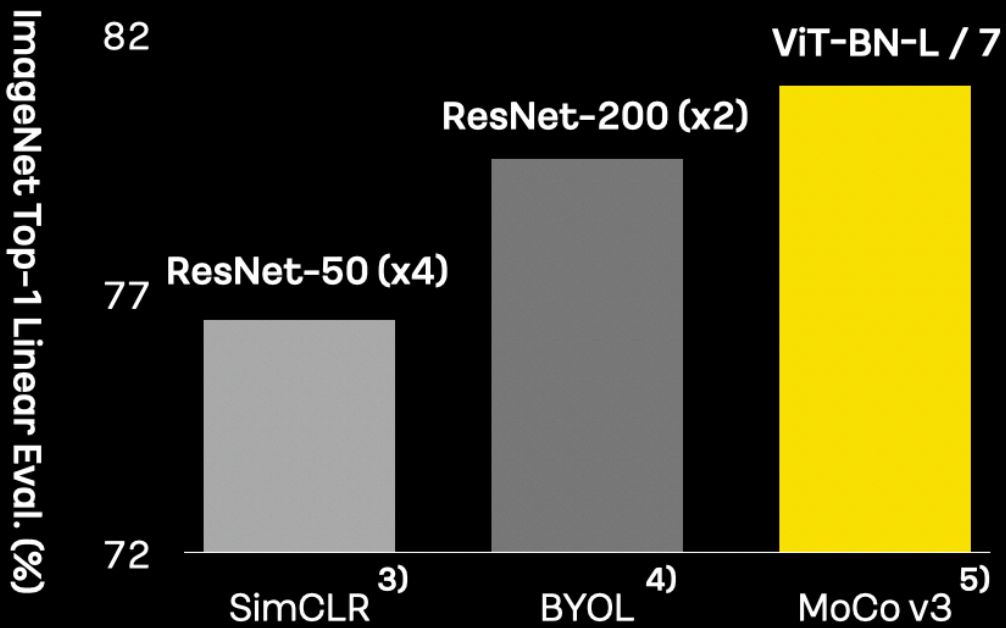How to automatically obtain similar and dissimilar pairs without labels?

# Contrastive Learning

Using a simple contrastive objective to learn global representations



$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\boldsymbol{z}_i, \boldsymbol{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\boldsymbol{z}_i, \boldsymbol{z}_k)/\tau)}$$

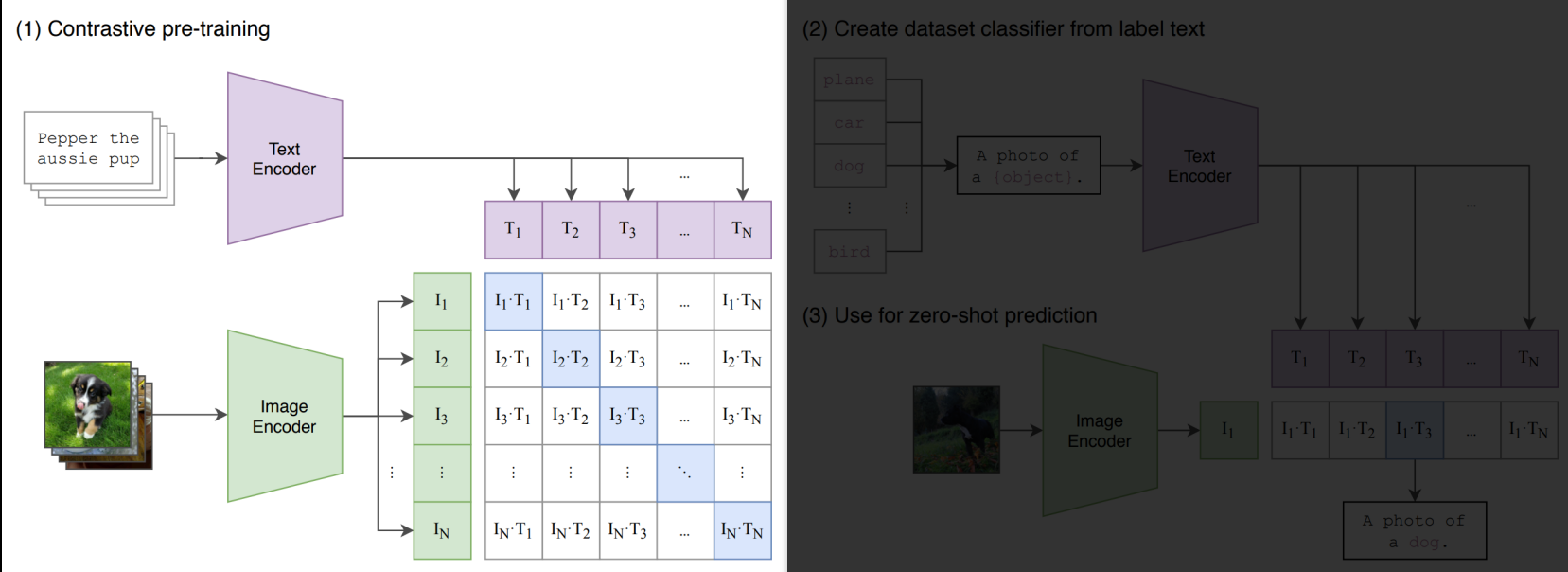# This simple approach really works well!

# What's the Next Step?

## Self-Supervised Multi-modal Representation Learning

# CLIP: Connecting Text and Images

Learning the shared global representations from images and texts!

Radford and Kim et al. "Learning Transferrable Visual Models from Natural Language Supervision", ICML'21.

# CLIP: Connecting Text and Images



(1) Contrastive pre-training

(2) Create dataset classifier from label text

(3) Use for zero-shot prediction

Radford and Kim et al. "Learning Transferrable Visual Models from Natural Language Supervision", ICML'21.

# and Images

```
# 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]
T_f = text_encoder(T)  #[n, d_t]

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# 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
```
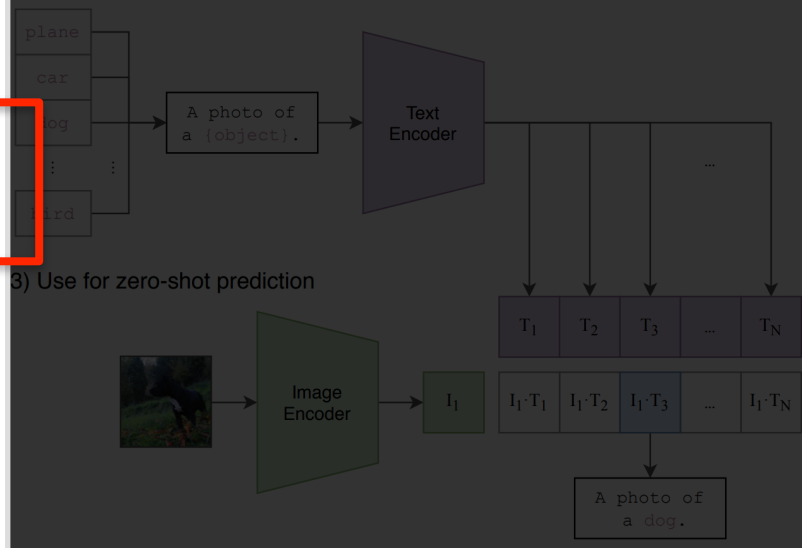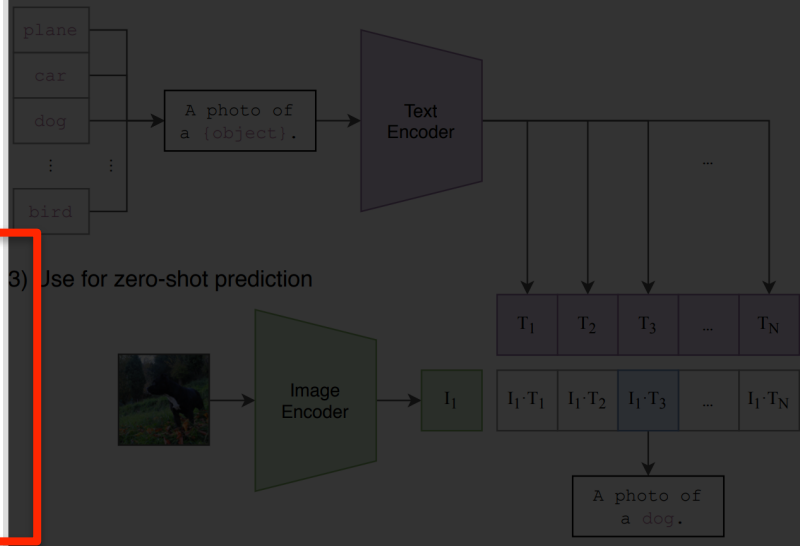
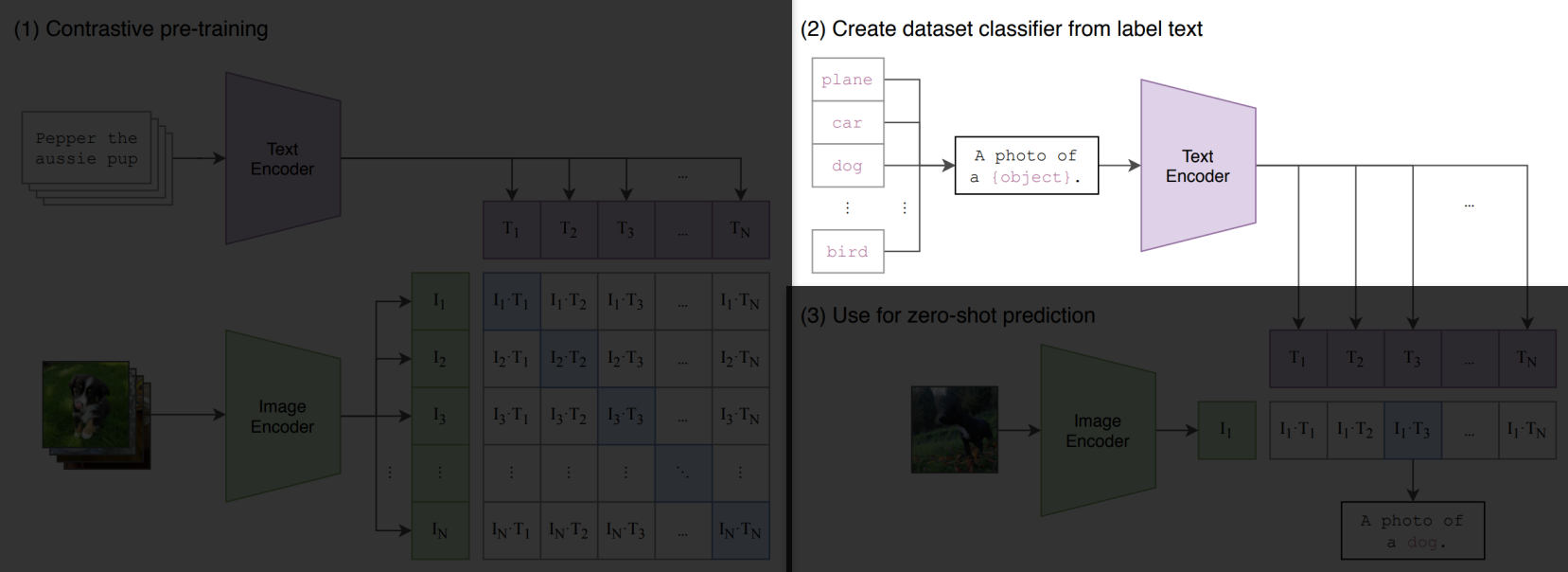*Figure 3.* Numpy-like pseudocode for the core of an implementation of CLIP.
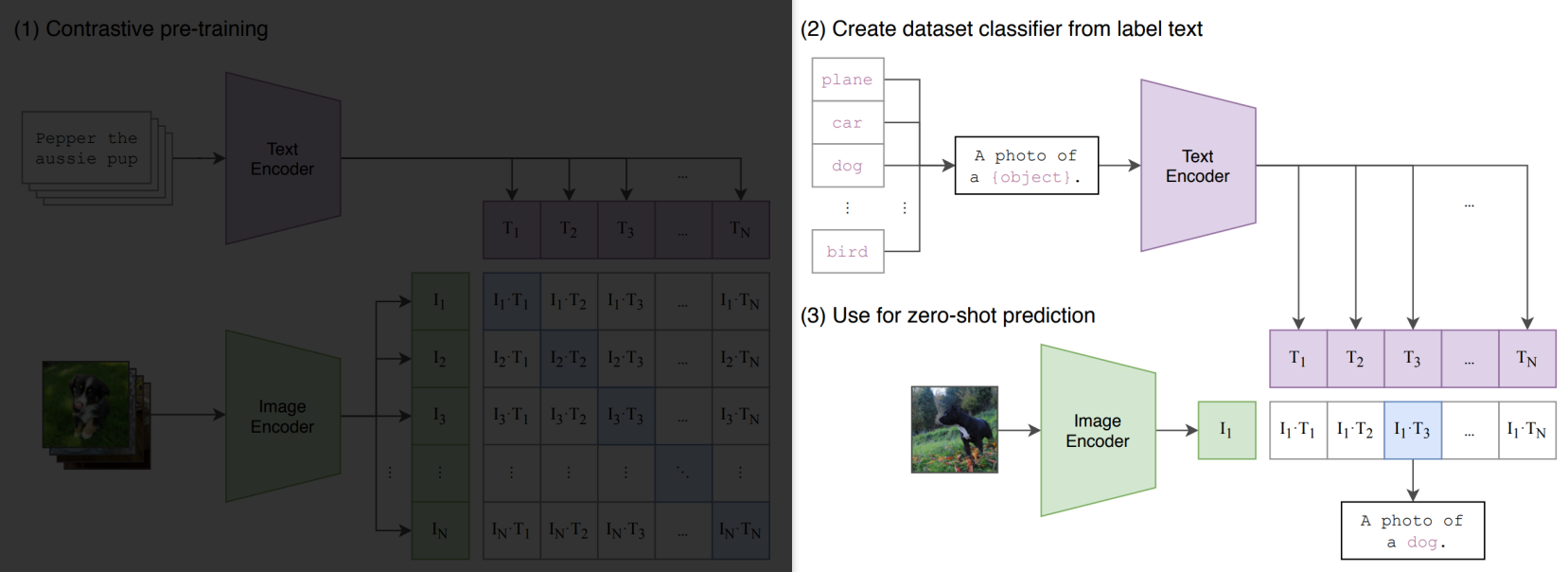


2) Create dataset classifier from label text

3) Use for zero-shot prediction

Radford and Kim et al. "Learning Transferrable Visual Models from Natural Language Supervision", ICML'21.

**and Images**

```
# 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]
T_f = text_encoder(T)  #[n, d_t]

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# 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
```

*Figure 3.* Numpy-like pseudocode for the core of an implementation of CLIP.



2) Create dataset classifier from label text

3) Use for zero-shot prediction

Radford and Kim et al. "Learning Transferrable Visual Models from Natural Language Supervision", ICML'21.

# and Images

```
# 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]
T_f = text_encoder(T)  #[n, d_t]

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# 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
```

*Figure 3.* Numpy-like pseudocode for the core of an implementation of CLIP.



2) Create dataset classifier from label text

3) Use for zero-shot prediction

Radford and Kim et al. "Learning Transferrable Visual Models from Natural Language Supervision", ICML'21.

# CLIP: Connecting Text and Images



(1) Contrastive pre-training

(2) Create dataset classifier from label text

(3) Use for zero-shot prediction

Radford and Kim et al. "Learning Transferrable Visual Models from Natural Language Supervision", ICML'21.

# CLIP: Connecting Text and Images



(1) Contrastive pre-training

(2) Create dataset classifier from label text

(3) Use for zero-shot prediction

Radford and Kim et al. "Learning Transferrable Visual Models from Natural Language Supervision", ICML'21.

# Large-scale Image-Text Pairs



MSCOCO sample

**"A shoe rack with some shoes and a dog sleeping on them"**

# Large-scale Image-Text Pairs

# Large-scale Image-Text Pairs

(CVPR'21)

# Large-scale Image-Text Pairs

# Large-scale Image-Text Pairs



Google Research (ICML'21)

# Where the Source?
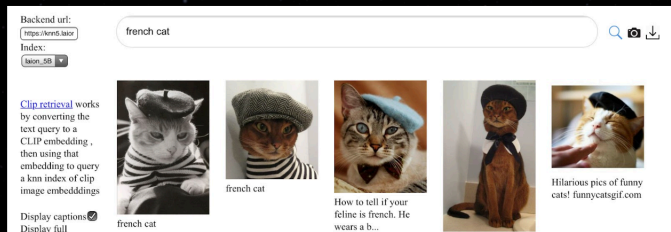
# LAION Projects

Romain Beaumont

31. March 2022

1 Comment
Uncategorized

## LAION-5B: A NEW ERA OF OPEN LARGE-SCALE MULTI-MODAL DATASETS

We present a dataset of 5,85 billion CLIP-filtered image-text pairs, 14x bigger than LAION-400M, previously the biggest openly accessible image-text dataset in the world.

Authors: Christoph Schuhmann, Richard Vencu, Romain Beaumont, Theo Coombes, Cade Gordon, Aarush Katta, Robert Kaczmarczyk, Jenia Jitsev

# Experiments — Zero-shot Classification



**+**

**A photo of {label}**

**A photo of a cat**

**A photo of a dog**

⋮

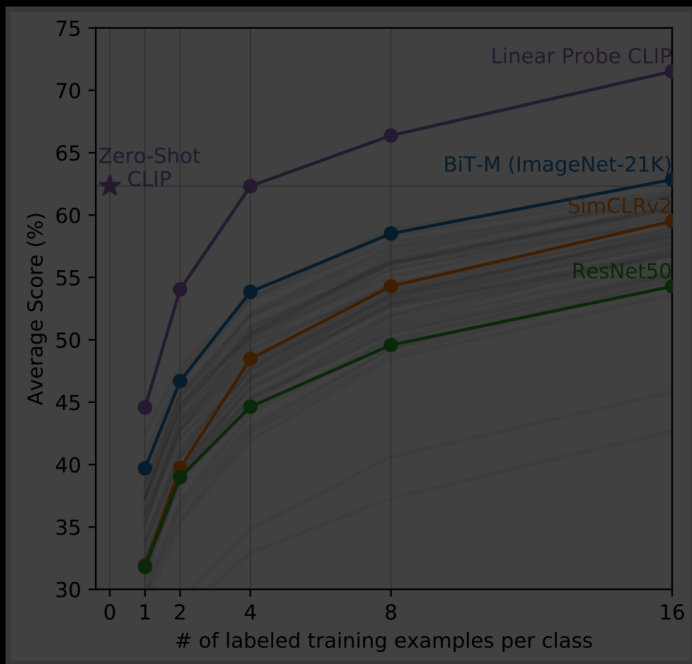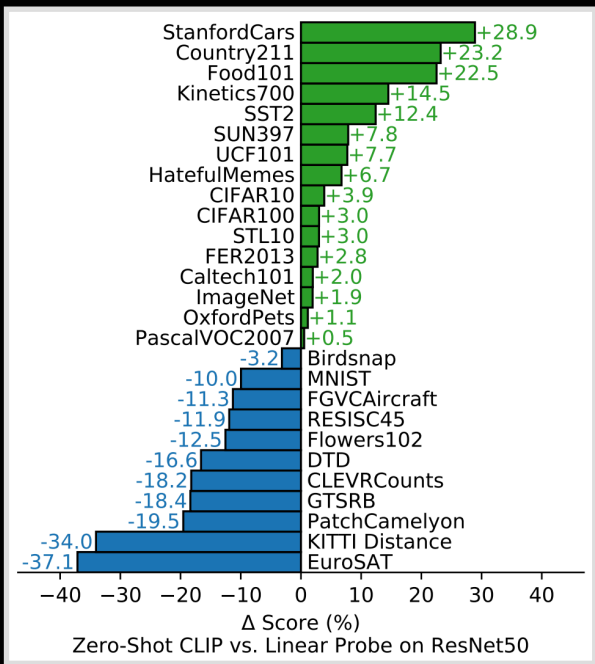# Experiments — Zero-shot Classification



**+**

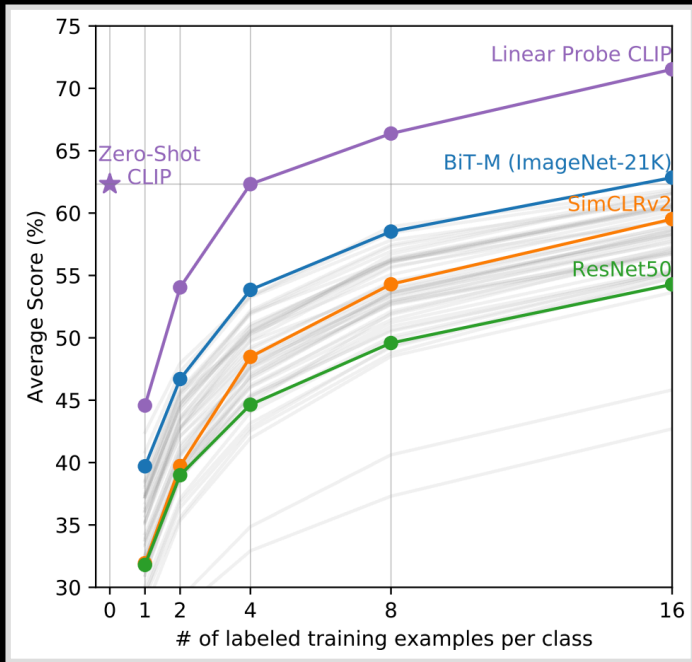A photo of {label}**, a type of flower**
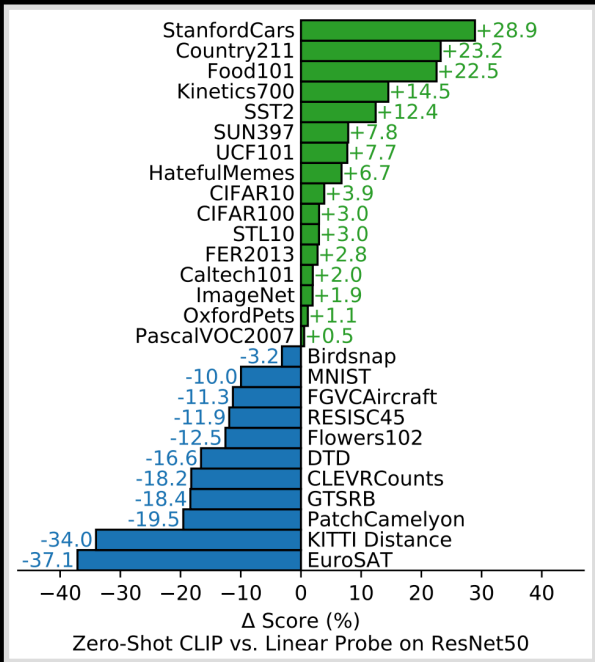
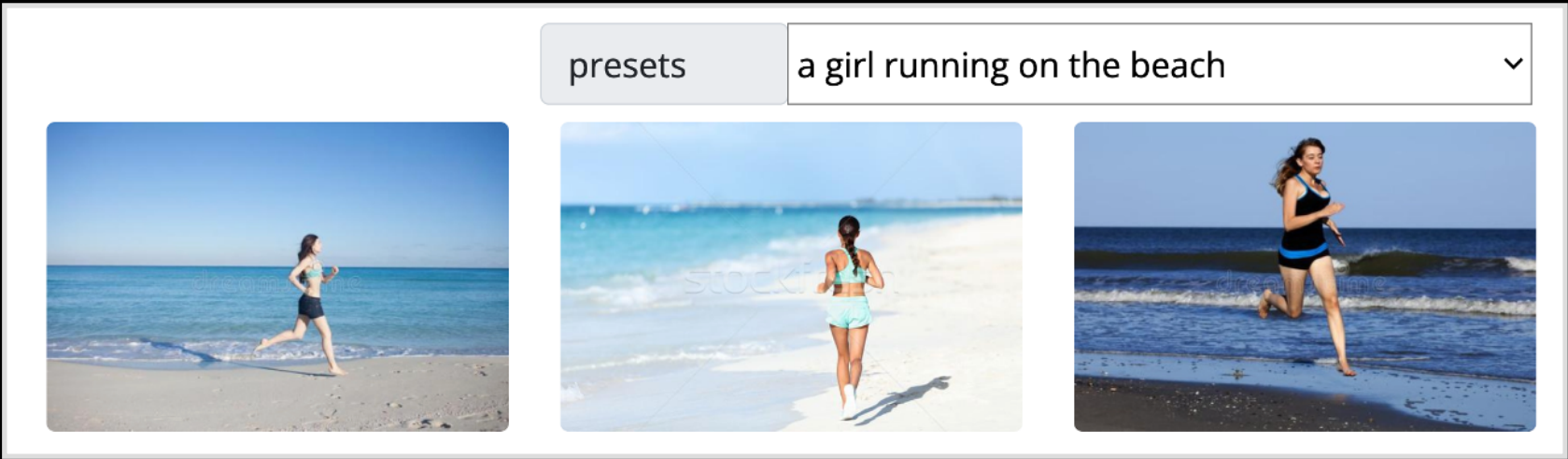A photo of a cat

A photo of a dog

⋮

# Experiments — Zero-shot Classification

# Experiments — Zero-shot Classification

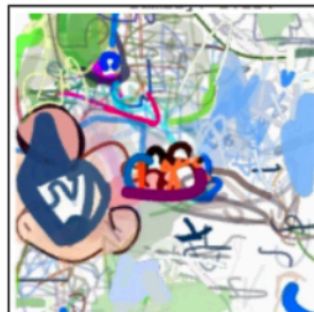# Application of CLIP — Search Engine

# Application of CLIP



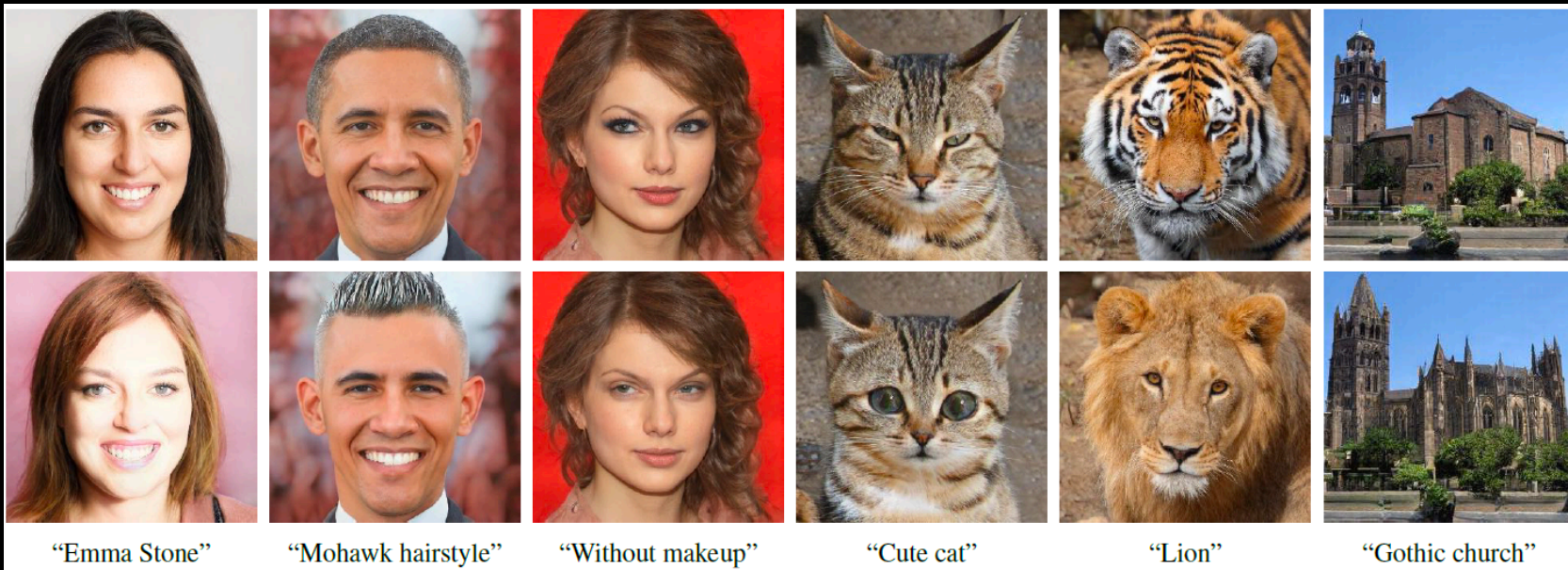"A drawing of a cat".   "Horse eating a cupcake".   "A 3D rendering of a temple".   "Family vacation to Walt Disney World".   "Self".

**Various drawings synthesized by CLIPDraw**, along with the corresponding description prompts used. CLIPDraw synthesizes images from text by performing gradient descent over a set of RGBA Bézier curves, with the goal of minimizing cosine distance between the CLIP encodings of generated images and description prompts. CLIPDraw does not require learning a new model, and can generally synthesize images within a minute on a typical GPU.

Frans et al. "CLIPDraw: Exploring Text-to-Drawing Synthesis through Language-Image Encoders", arXiv'21

# Application of CLIP



"Emma Stone"     "Mohawk hairstyle"     "Without makeup"     "Cute cat"     "Lion"     "Gothic church"

# Conclusion

Now, it's possible to learn a shared representation from text-image pairs