



Learning Transferable Visual Models From Natural Language Supervision

Main presentation

Alec Radford, Jong Wook Kim, Chris Hallacy, A. Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, Ilya Sutskever, *OpenAI*, 2021.

Presenter:

Diego Calanzone

Seminar:

Learning with Limited Labeled Data

Academic year:

2021/2022



Overview

- **Natural Language Supervision: SoTA, ZS performance**

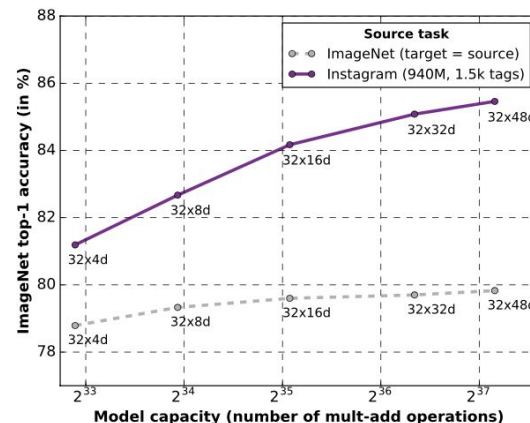
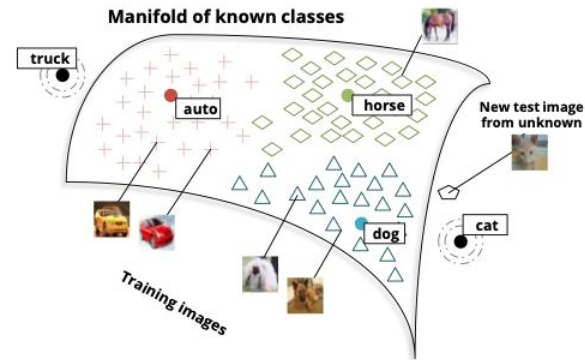
- Mapping images to text embeddings
- Learning better image representations
- Exploiting large web data

- **CLIP : motivation, intuition**

- Dataset & data pipeline
- Experimental setup, methodological components
- Scaling, data efficiency
- Performance, robustness, bias
- Limitations and broader impacts

- **CLIP “spin-offs”:**

- Multimodal neurons
- Zero-shot text generation: DALL-E

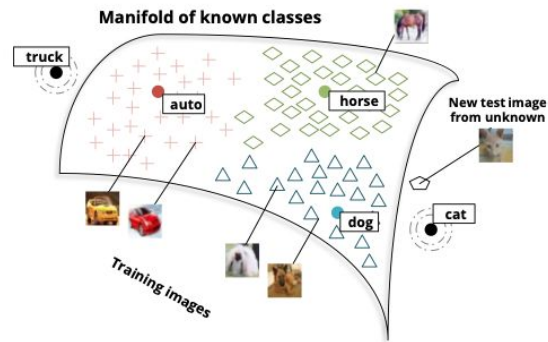




Natural Language Supervision (Supervised)

Zero-Shot Learning Through Cross-Modal Transfer

R. Socher, M. Ganjoo, C. D. Manning, A. Y. Ng, 2013



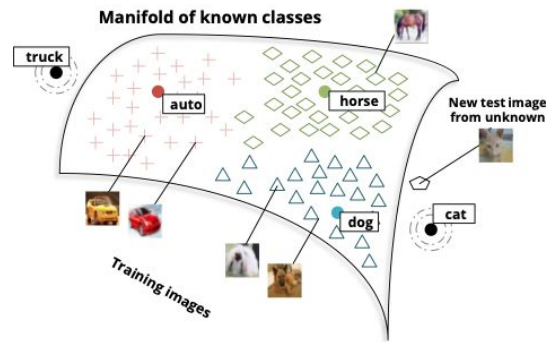


Natural Language Supervision (Supervised)

Zero-Shot Learning Through Cross-Modal Transfer

R. Socher, M. Ganjoo, C. D. Manning, A. Y. Ng, 2013

- Training dataset: CIFAR-100
- Key points: semantic embedding space, novelty detection
- ZS Class. Accuracy: **52.7%** (CIFAR-100, **6** novel classes)



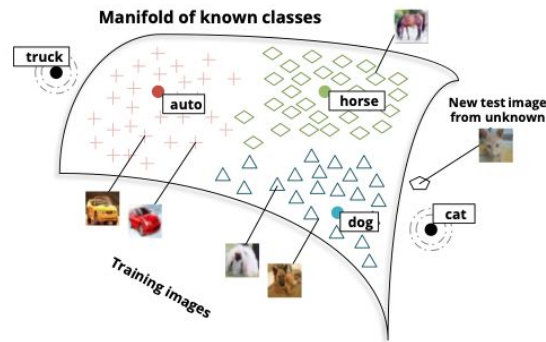


Natural Language Supervision (Supervised)

Zero-Shot Learning Through Cross-Modal Transfer

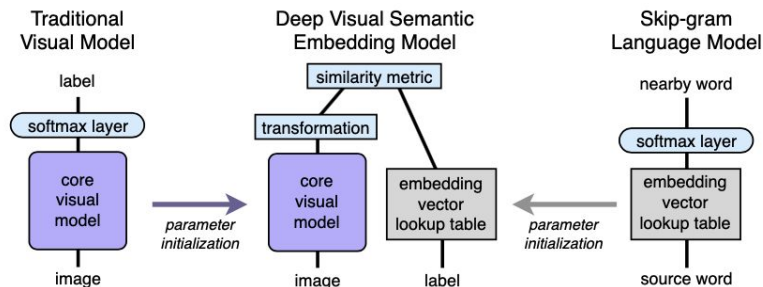
R. Socher, M. Ganjoo, C. D. Manning, A. Y. Ng, 2013

- Training dataset: CIFAR-100
- Key points: semantic embedding space, novelty detection
- ZS Class. Accuracy: **52.7%** (CIFAR-100, 6 novel classes)



DeViSE: A Deep Visual-Semantic Embedding Model

A. Frome, G. S. Corrado*, J. Shlens, et al., 2013



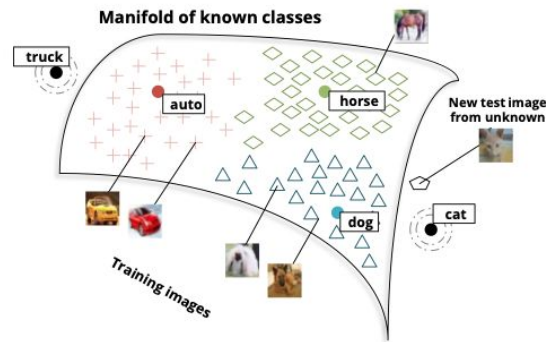


Natural Language Supervision (Supervised)

Zero-Shot Learning Through Cross-Modal Transfer

R. Socher, M. Ganjoo, C. D. Manning, A. Y. Ng, 2013

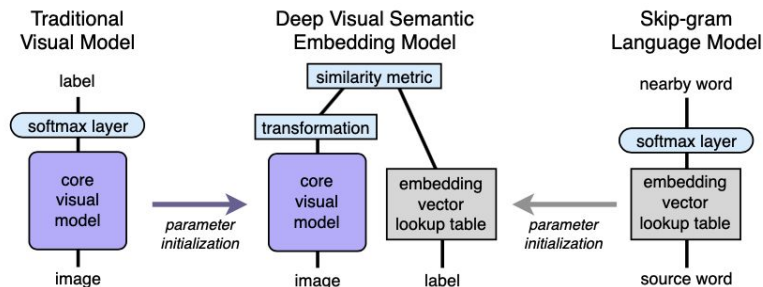
- Training dataset: CIFAR-100
- Key points: semantic embedding space, novelty detection
- ZS Class. Accuracy: **52.7%** (CIFAR-100, 6 novel classes)



DeViSE: A Deep Visual-Semantic Embedding Model

A. Frome, G. S. Corrado*, J. Shlens, et al., 2013

- Training dataset: ILSVRC 2012 1K
- Key points: semantic vector prediction, language-model supervision
- ZS Class. Accuracy: **36.4%**, ImageNet 2011 21K

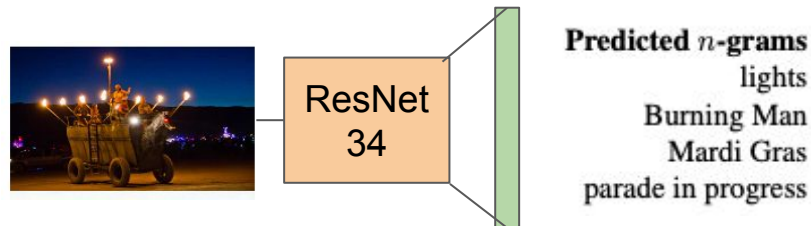




Natural Language Supervision (Supervised)

Learning Visual N-Grams from Web Data

A. Li, A. Jabri, A. Joulin, L. van der Maaten, 2017



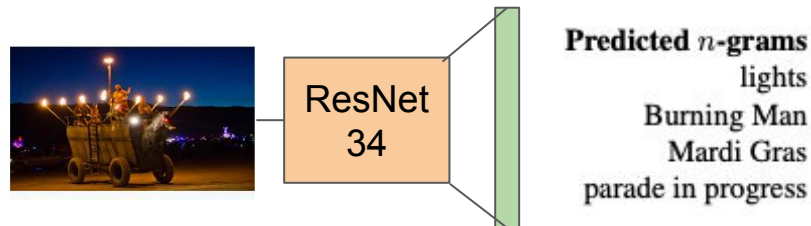


Natural Language Supervision (Supervised)

Learning Visual N-Grams from Web Data

A. Li, A. Jabri, A. Joulin, L. van der Maaten, 2017

- Training dataset: YFCC100M (image-comments)
- Key points: ImageNet fine-tuned, n-grams probability distribution, conditional probability
- ZS Classif. Accuracy: **11.5%** (ImageNet-1K), **23.0%** (SUN), **72.4%** (Yahoo)



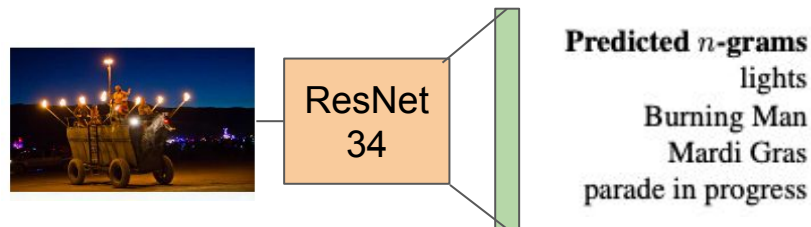


Natural Language Supervision (Supervised)

Learning Visual N-Grams from Web Data

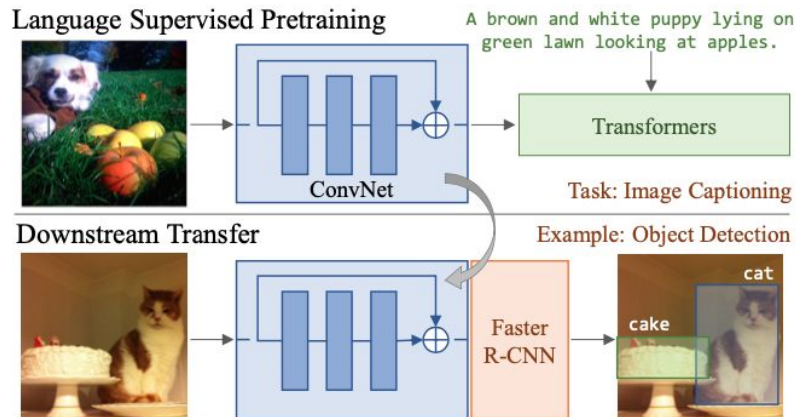
A. Li, A. Jabri, A. Joulin, L. van der Maaten, 2017

- Training dataset: YFCC100M (image-comments)
- Key points: ImageNet fine-tuned, n-grams probability distribution, conditional probability
- ZS Classif. Accuracy: **11.5%** (ImageNet-1K), **23.0%** (SUN), **72.4%** (Yahoo)



VirTex: Learning Visual Repr. from Text Annotations

K. Desai, J. Johnson, 2020



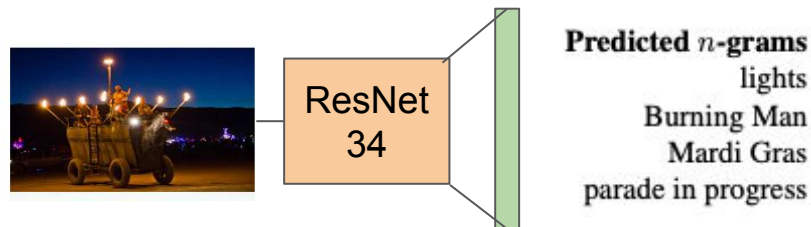


Natural Language Supervision (Supervised)

Learning Visual N-Grams from Web Data

A. Li, A. Jabri, A. Joulin, L. van der Maaten, 2017

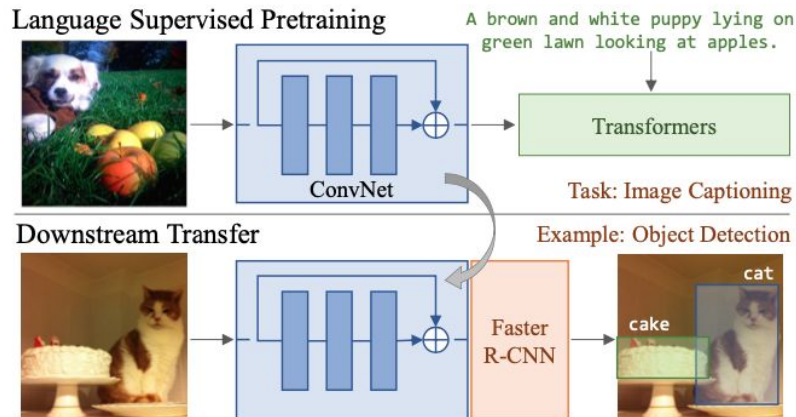
- Training dataset: YFCC100M (image-comments)
- Key points: ImageNet fine-tuned, n-grams probability distribution, conditional probability
- ZS Classif. Accuracy: **11.5%** (ImageNet-1K), **23.0%** (SUN), **72.4%** (Yahoo)



VirTex: Learning Visual Repr. from Text Annotations

K. Desai, J. Johnson, 2020

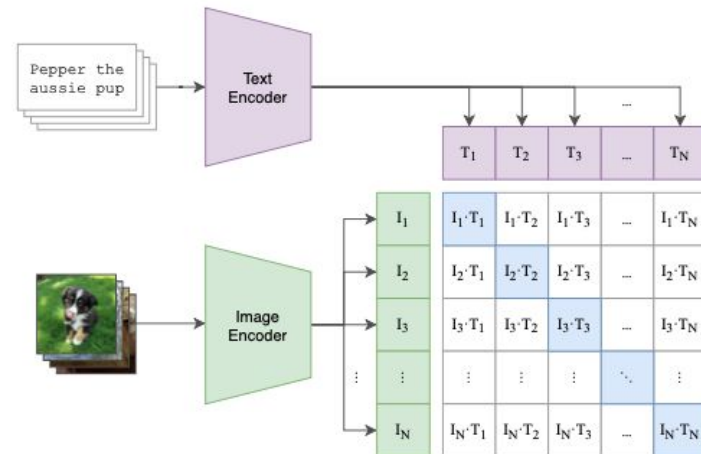
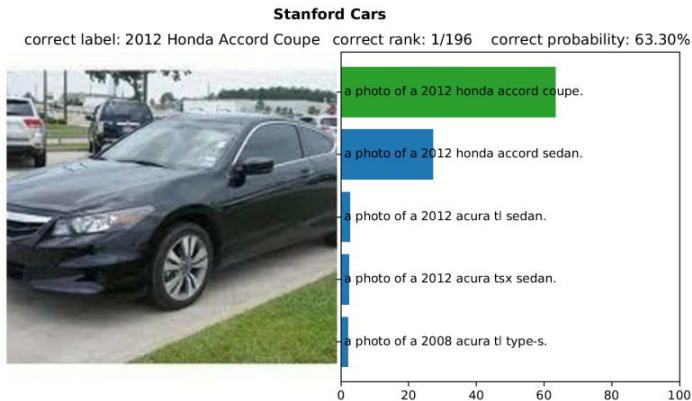
- Training dataset: COCO-captions (train2017 split)
- Key points: transformer cross-attention with CNN encodings
- Classif. Accuracy: **88.7%** (PASCAL VOC07), **53.8%** (ImageNet-1K)





Overview

- **Natural Language Supervision: SoTA, ZS performance**
 - Mapping images to text embeddings
 - Learning better image representations
 - Exploiting large web data
- **CLIP : motivation, intuition**
 - Dataset & data pipeline
 - Experimental setup, methodological components
 - Scaling, data efficiency
 - Performance, robustness, bias
 - Limitations and broader impacts
- **CLIP “spin-offs”:**
 - Multimodal neurons
 - Zero-shot text generation: DALL-E





Choice of data & pre-processing pipeline

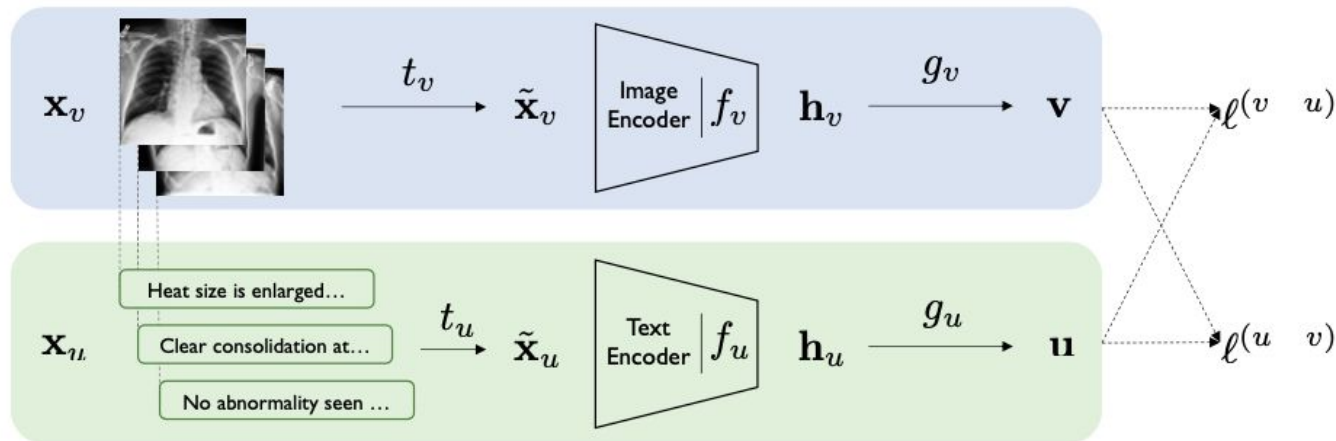
- **Objective:** high performance in zero-shot and new downstream tasks
MS-COCO, Visual Genome → too **small**; YFCC100M → **arguable** quality
- **Choice:** creating a new dataset **“WIT: WebImageText” 400M**
 - (image, text) pairs from a list of **500k queries**: words occurring >100 times on Wikipedia
 - Limit of **20k (img-txt) pairs / query**
- **Prompt engineering + Ensembling:** tackling polysemy and averaging multiple prompts
ie: “a boxer, a type of dog” vs. “a boxer, a type of athlete”
Accuracy gain → + ~5%
- **Data overlap analysis:** threshold on image distances in the embedding space → duplicate detection → 3.2% avg. overlap → +0.6% accuracy gain



The foundations of CLIP

Contrastive Learning of Medical Visual Representations from Paired Images and Text

Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D. Manning, Curtis P. Langlotz, 2020





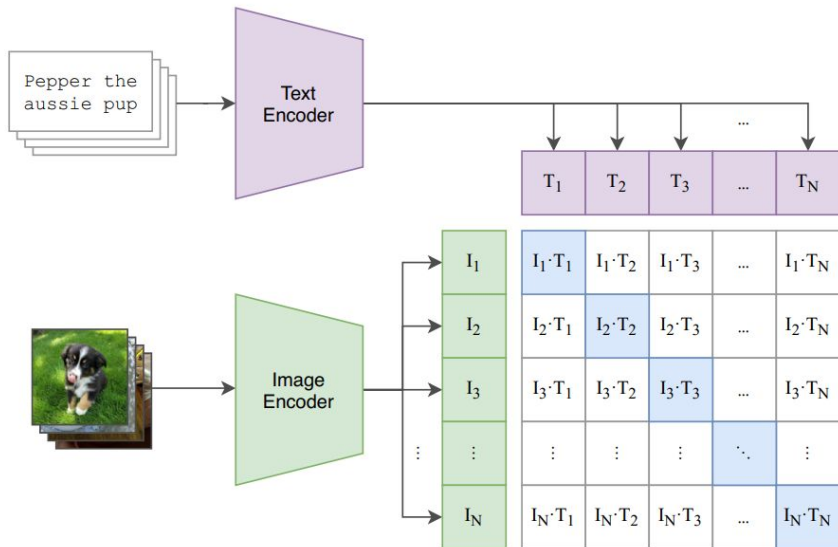
The foundations of CLIP

- Image/text embedding vectors: $\mathbf{v} = g_v(f_v(\tilde{\mathbf{x}}_v)) \quad \mathbf{u} = g_u(f_u(\tilde{\mathbf{x}}_u))$
- Image \rightarrow text contrastive loss: $\ell_i^{(v \rightarrow u)} = -\log \frac{\exp(\langle \mathbf{v}_i, \mathbf{u}_i \rangle / \tau)}{\sum_{k=1}^N \exp(\langle \mathbf{v}_i, \mathbf{u}_k \rangle / \tau)}$
- Text \rightarrow image contrastive loss: $\ell_i^{(u \rightarrow v)} = -\log \frac{\exp(\langle \mathbf{u}_i, \mathbf{v}_i \rangle / \tau)}{\sum_{k=1}^N \exp(\langle \mathbf{u}_i, \mathbf{v}_k \rangle / \tau)}$
- Loss function: $\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \left(\lambda \ell_i^{(v \rightarrow u)} + (1 - \lambda) \ell_i^{(u \rightarrow v)} \right)$

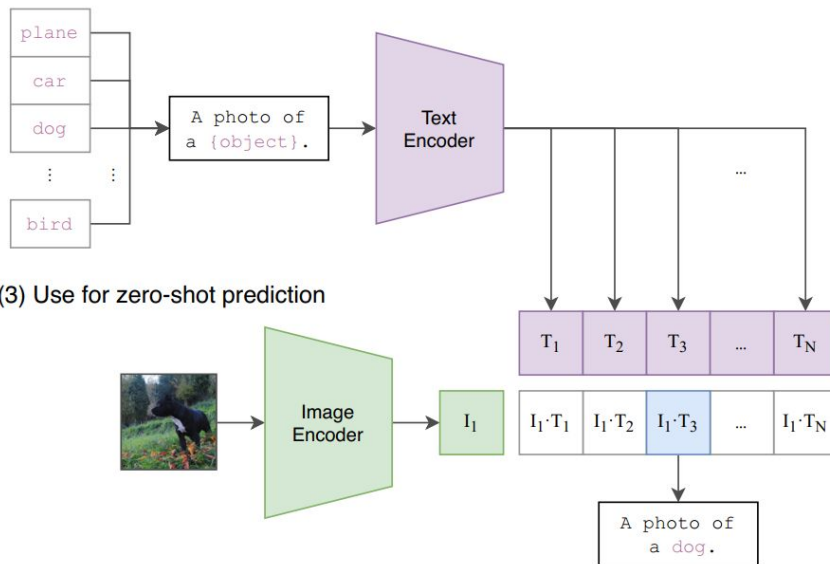


Contrastive Language-Image Pretraining (CLIP)

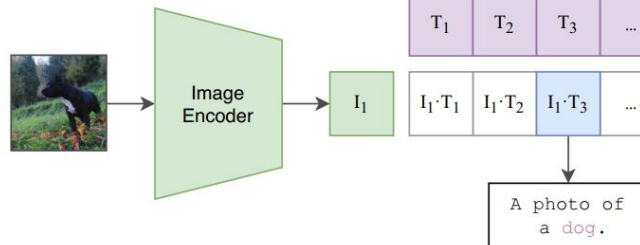
(1) Contrastive pre-training



(2) Create dataset classifier from label text

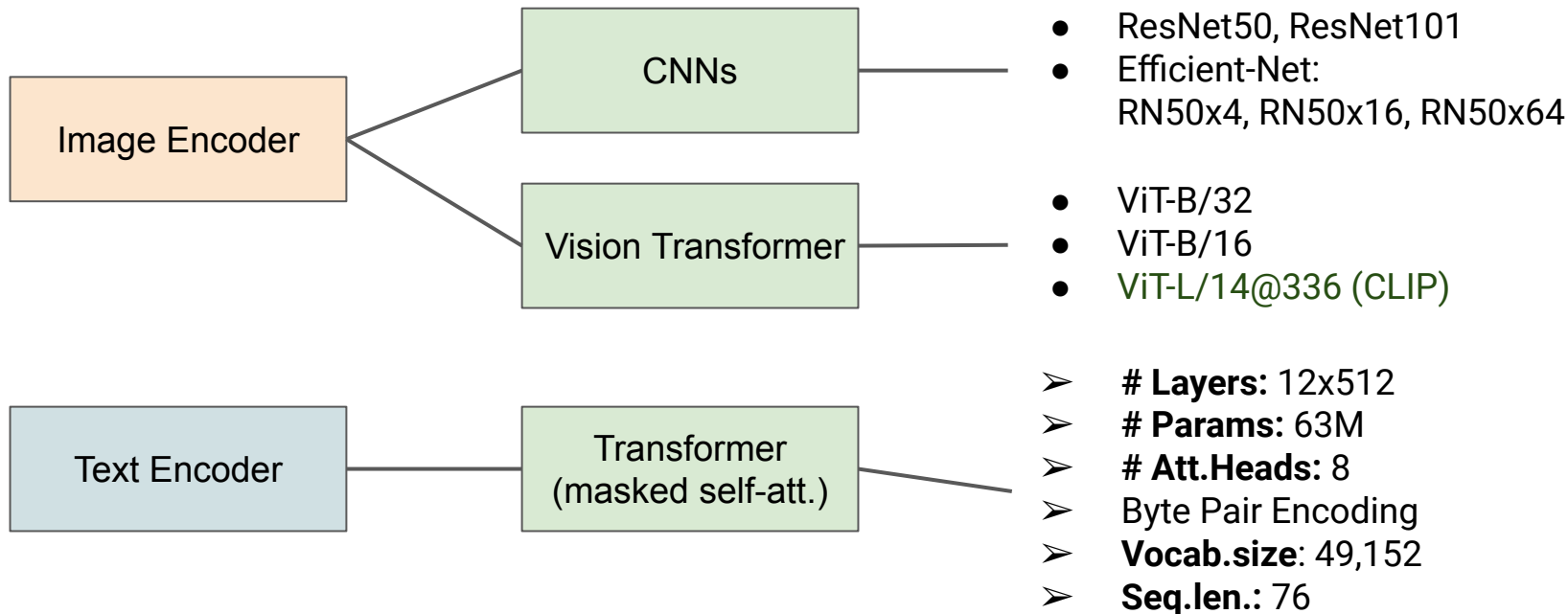


(3) Use for zero-shot prediction





Training setup & choice of the architecture

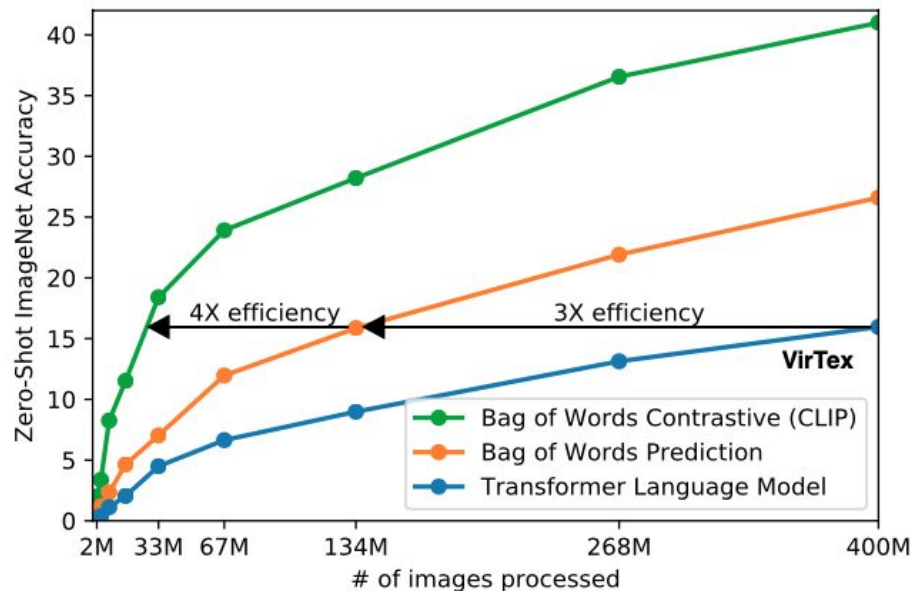


- **Optimizer:** Adam+reg., **mini-batch size:** 32,768, **epochs:** 32



Efficiency in scaling the dataset

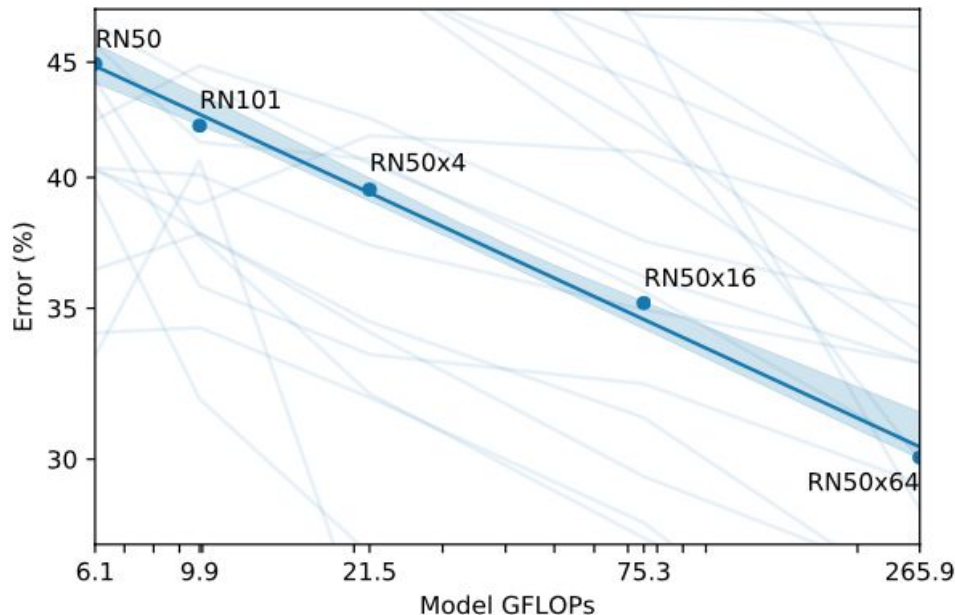
- **Data scaling:** CLIP > VirTex, bag of words > sentences





Efficiency in scaling the model

- **Data scaling:** CLIP > VirTex, bag of words > sentences
- **Model scaling:** up to 44x → smooth error decrease (like GPT). Scaling:
 - img: width, depth, res.
 - txt: width

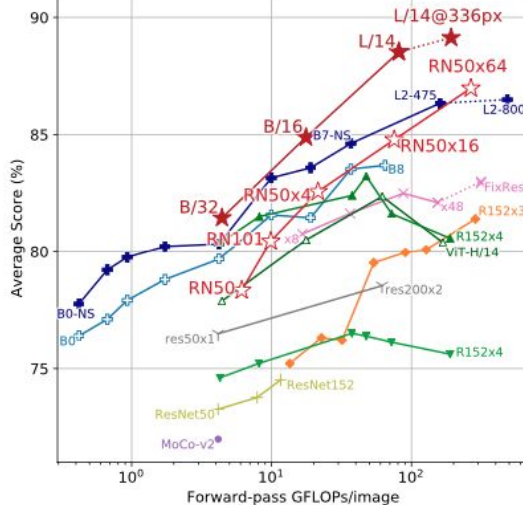




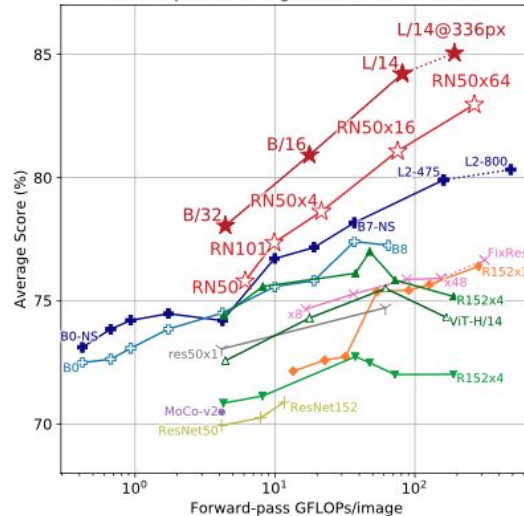
Compared classification performance in scaling

- **Data scaling:** CLIP > ViT, bag of words > sentences
- **Model scaling:** up to 44x → smooth error decrease (like GPT). Scaling:
 - img: width, depth, res.
 - txt: width
- Linear probe
CLIP-ViT-L/14-336px
outperforms the CNN-baseline

Linear probe average over Kornblith et al.'s 12 datasets



Linear probe average over all 27 datasets

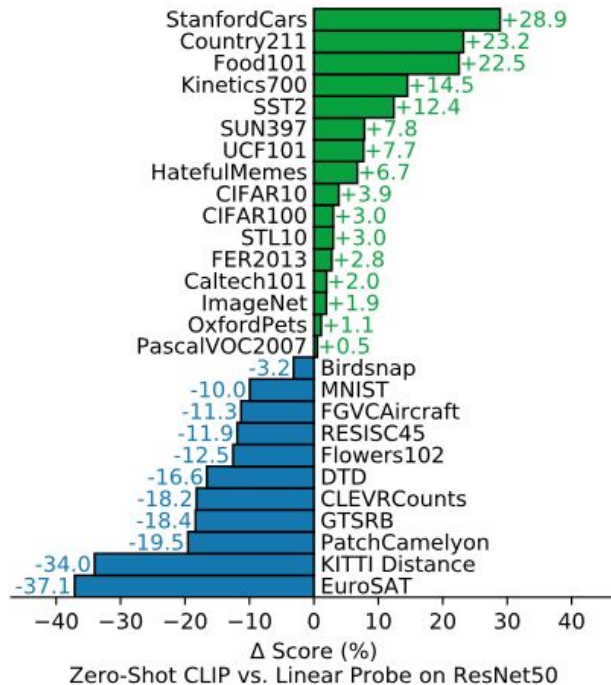




Zero-shot Performance

- StanfordCars, STL10: SoTA performance on datasets with few label examples and **a lot of unlabeled data**.
- Kinetics700, UCF101: NLS → better context with verbs
- EuroSAT, PatchCamelyon: fine-grained tasks, specific examples → weak

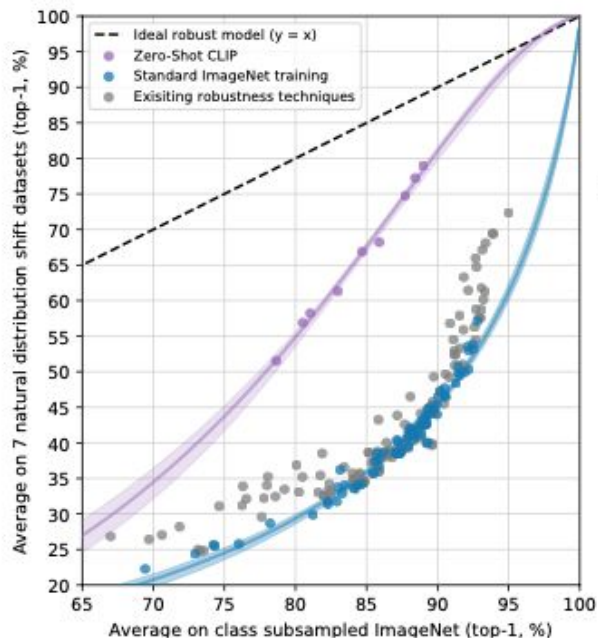
Advantages with CLIP : output dynamicity, better text-supervision, pre-trained once.



	aYahoo	ImageNet	SUN
Visual N-Grams	72.4	11.5	23.0
CLIP	98.4	76.2	58.5



Robustness with distributional shifts



	Dataset Examples						ImageNet ResNet101	Zero-Shot CLIP	Δ Score
ImageNet							76.2	76.2	0%
ImageNetV2							64.3	70.1	+5.8%
ImageNet-R							37.7	88.9	+51.2%
ObjectNet							32.6	72.3	+39.7%
ImageNet Sketch							25.2	60.2	+35.0%
ImageNet-A							2.7	77.1	+74.4%



Limitations & Impacts

- Required **1000x** scaling → SOTA



Limitations & Impacts

- Required **1000x** scaling → SOTA
- Poor performance in specific/more **abstract** tasks



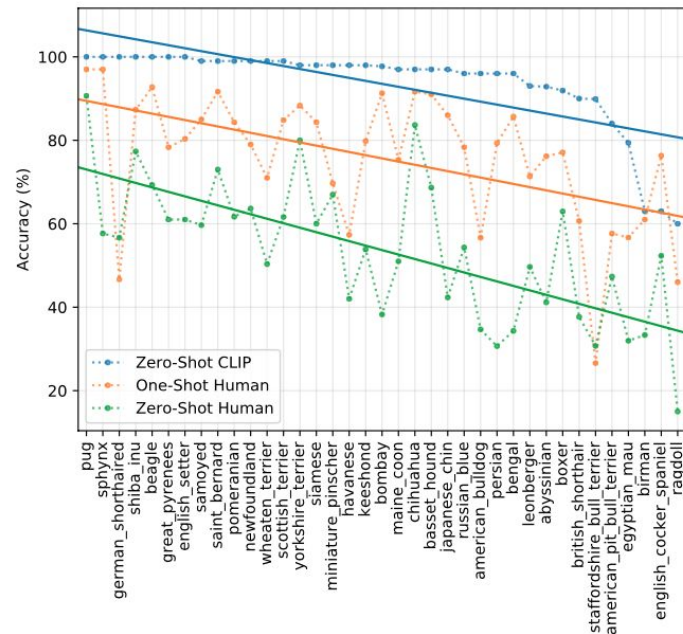
Limitations & Impacts

- Required **1000x** scaling → SOTA
- Poor performance in specific/more **abstract** tasks
- Underperforms with **synthetic** data → only 88% ZS acc. on MNIST!



Limitations & Impacts

- Required **1000x** scaling → SOTA
- Poor performance in specific/more **abstract** tasks
- Underperforms with **synthetic** data → only 88% ZS acc. on MNIST!
- Unfiltered dataset → **bias** harm
- Inefficient few-shot performance wrt. humans



Broader Impacts

- “Roll your own classifier” without re-training
- Surveillance: face recognition, emotion recognition, action recognition, geo-localization

	Accuracy	Majority Vote on Full Dataset	Accuracy on Guesses	Majority Vote Accuracy on Guesses
Zero-shot human	53.7	57.0	69.7	63.9
Zero-shot CLIP	93.5	93.5	93.5	93.5
One-shot human	75.7	80.3	78.5	81.2
Two-shot human	75.7	85.0	79.2	86.1



“Spin-offs” of CLIP

- CLIP for Surveillance
 - Search through frames in video
 - “Roll your own classifier”
- Multi-modal neurons found in CLIP
 - Neurons that respond to **particular inputs**
 - Work **similarly** to the **brain’s** neurons
 - They organize **highly abstract concepts**
- Zero-shot image generation: DALL-E
 - VQ-VAE + GPT-style Transformer
 - CLIP to rank generated images

Christmas



Any



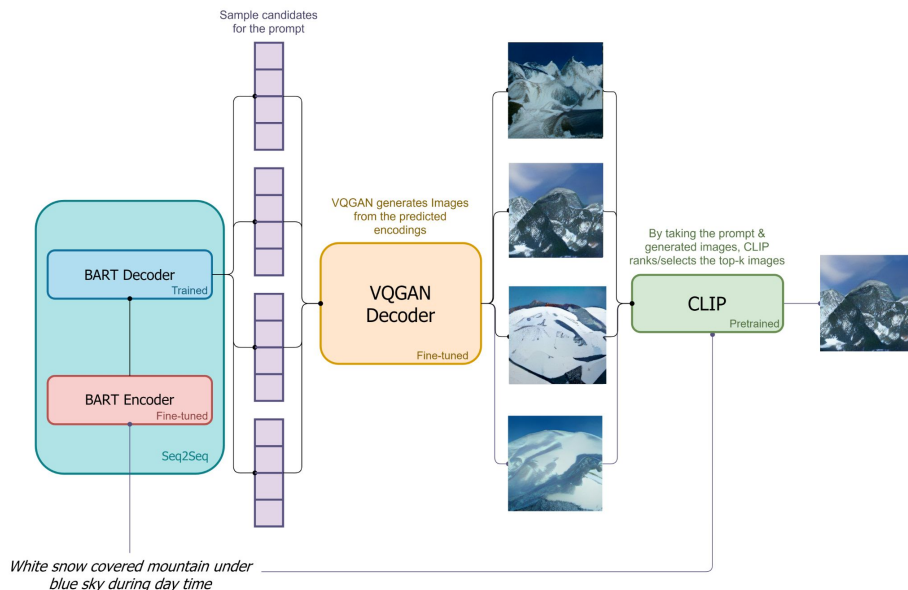
Text



Face



Logo



Thanks for listening!

Questions are welcome.

*(More on DALL-E & multi-modal neurons
in the Q&A if requested)*

TEXT PROMPT

a store front that has the word 'openai' written on it. ...

AI-GENERATED IMAGES



West Africa



Any

Text

Face

Logo

USA



Any

Text

Face

Logo



For notebook examples, further insights, sources, **I made a repo!**

<https://github.com/halixness/understanding-CLIP>



Multimodal neurons in CLIP : visualizing concepts

BONUS

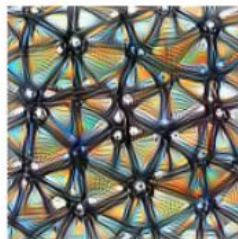
By **optimization** (ie. with GANs/ AEs) → find images that **maximize the activation** of a component

<https://distill.pub/2017/feature-visualization/>



Neuron

$\text{layer}_n[x, y, z]$



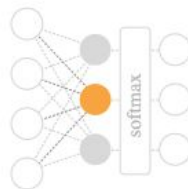
Channel

$\text{layer}_n[:, :, z]$



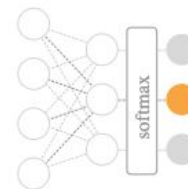
Layer/DeepDream

$\text{layer}_n[:, :, :]^2$



Class Logits

$\text{pre_softmax}[k]$



Class Probability

$\text{softmax}[k]$



Multimodal neurons in CLIP

BONUS

heart



Any



Text



Face



Logo



Architecture



Indoor



Nature



Pose



Multimodal neurons in CLIP

BONUS

USA



Any



Text



Face



Logo



Architecture



Indoor



Nature

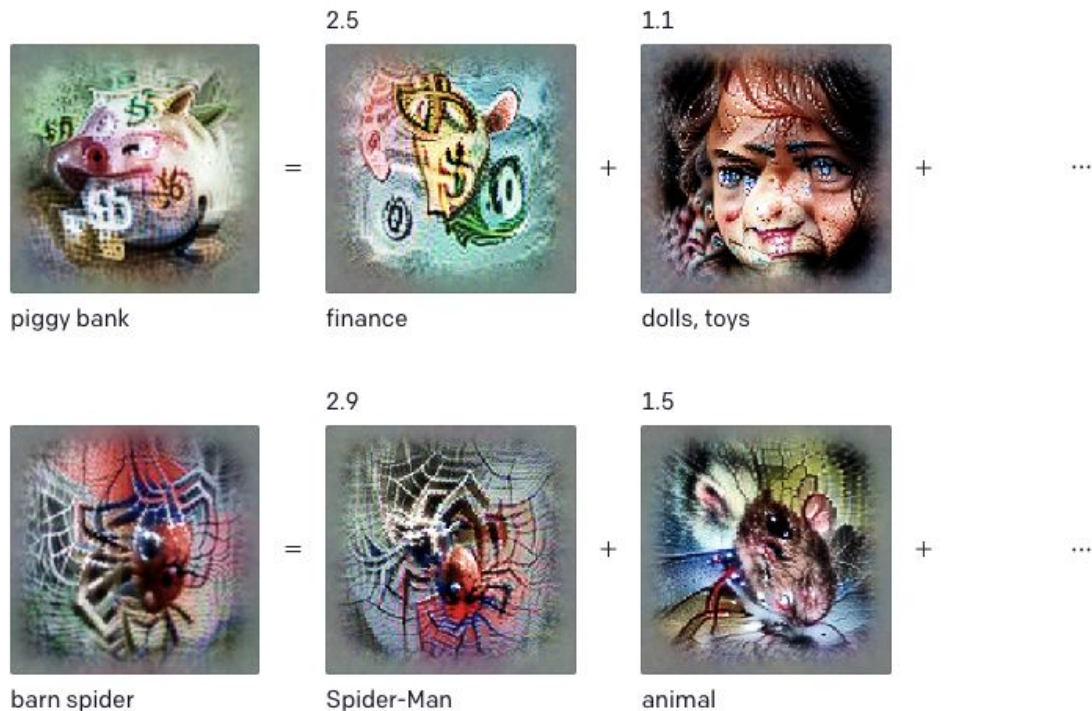


Pose



Multimodal neurons in CLIP

BONUS

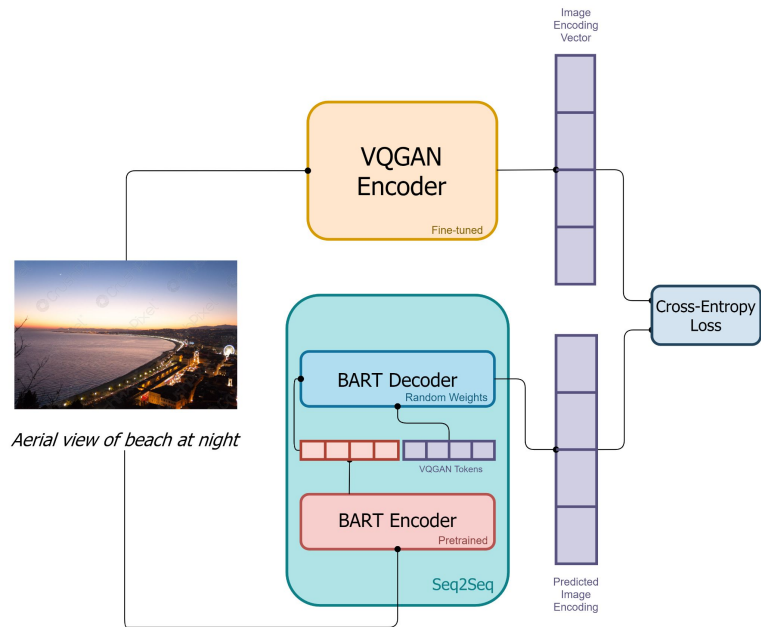




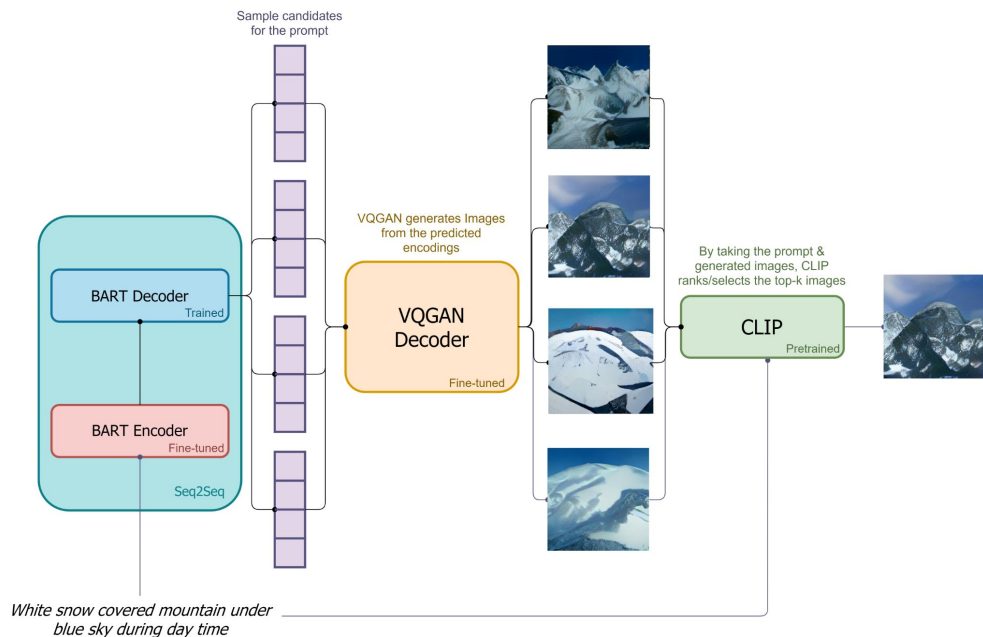
Zero-shot text generation: DALL-E

BONUS

Training



Inference





Zero-shot text generation: DALL-E

BONUS

a farmhouse surrounded by beautiful flowers



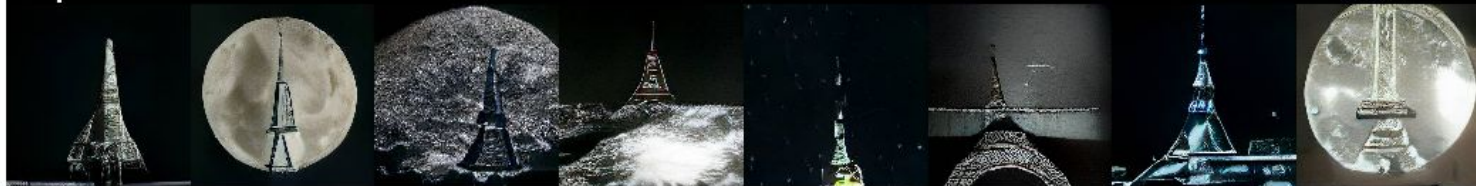
pred-05

sunset over green mountains



pred-07

a picture of the Eiffel tower on the Moon





Zero-shot text generation: DALL-E

BONUS

TEXT PROMPT

an armchair in the shape of an avocado. ...

AI-GENERATED IMAGES

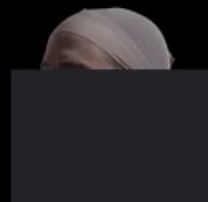


Edit prompt or view more images ↗

TEXT PROMPT

a photograph of a bust of homer

IMAGE PROMPT



AI-GENERATED IMAGES

