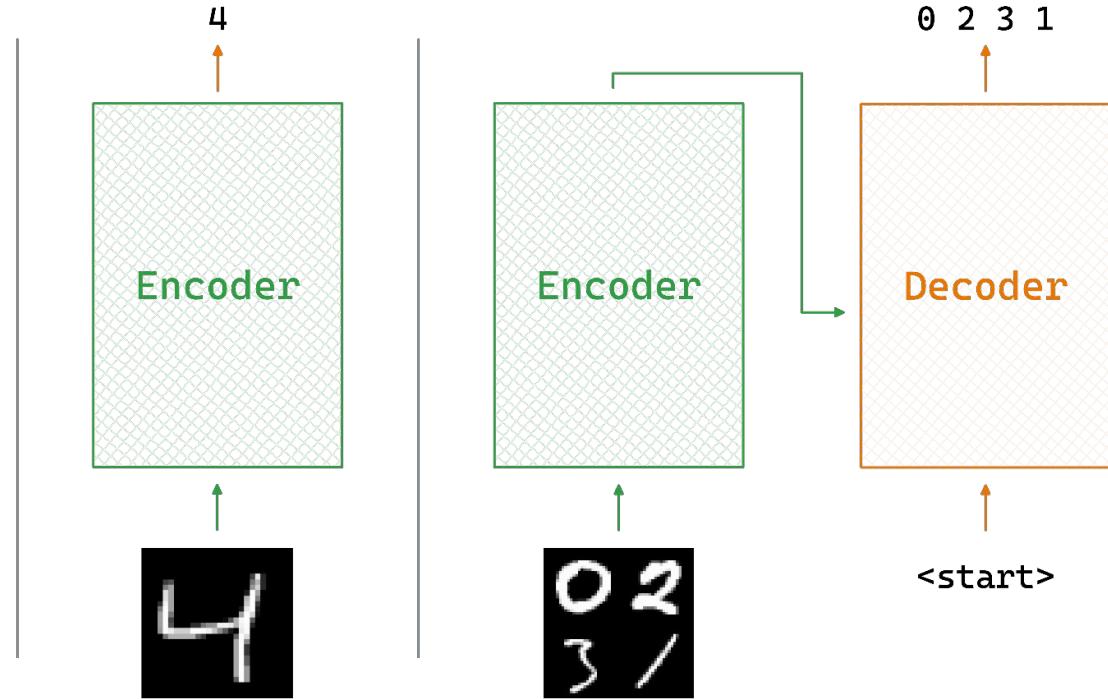


Multimodal transfer learning

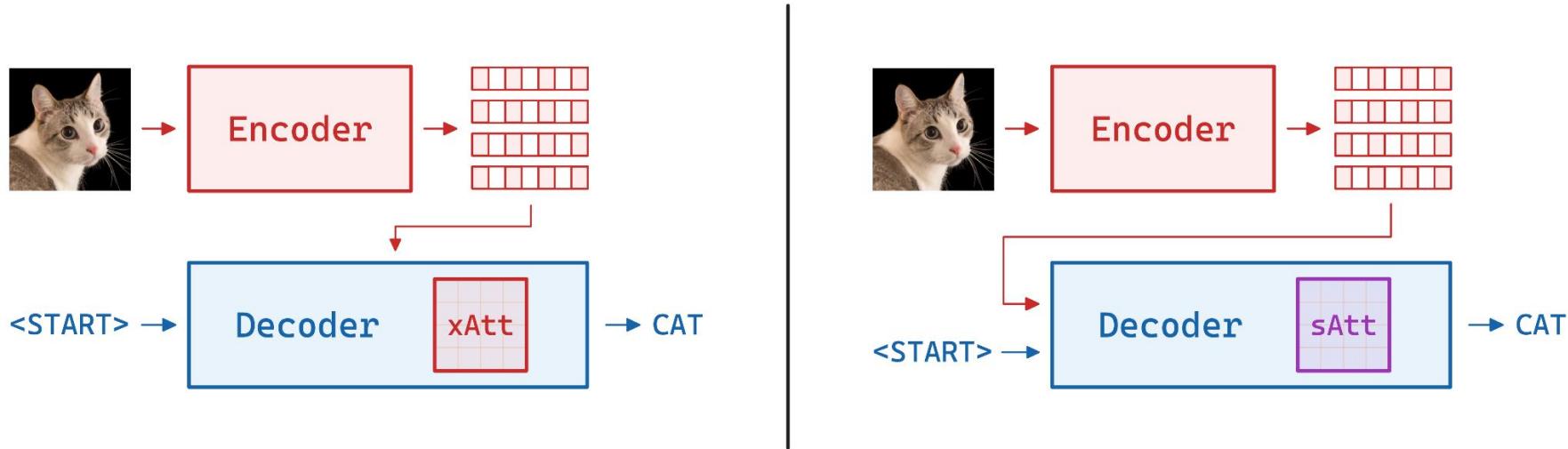
Merge models for “novel” tasks

Previous Task

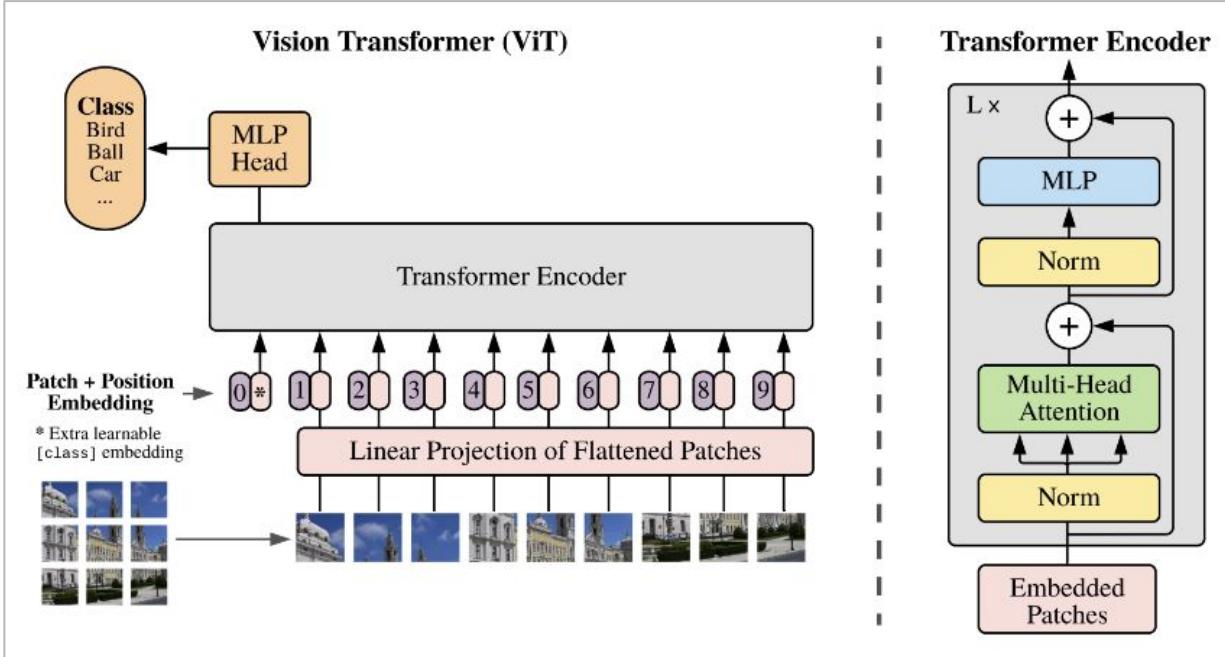
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9



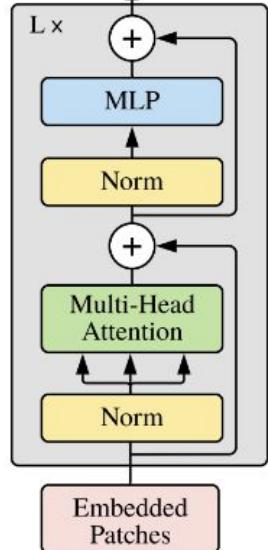
Transfer Learning



Vision Transformer (ViT)



Transformer Encoder



Published as a conference paper at ICLR 2021

AN IMAGE IS WORTH 16x16 WORDS:
TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy¹, Lucas Beyer², Alexander Kolesnikov³, Dirk Weissenborn², Xiaohua Zhai², Thomas Unterthiner², Matthias Minderer², Matthias Minderer², Georg Heigold², Sylvain Gelly², Jakob Uszkoreit², Neil Houlsby²
¹Equal technical contribution
²Google Research, Brain Team
{adosovitskiy, neilhoulsby}@google.com

ABSTRACT

While the Transformer architecture has become the de-facto standard for natural language processing (NLP), its application to computer vision tasks is limited. In vision, attention is either used in conjunction with convolutional networks, or used to replace certain components of convolutional networks while keeping their visual structure intact. We show that the multi-head self-attention mechanism and a linear transformer applied directly to sequences of image patches can perform very well on image classification tasks. When pre-trained on large amounts of data and fine-tuned on small datasets, our model outperforms state-of-the-art baselines (ImageNet, CIFAR-10, VTFB, etc.). Our Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring substantially fewer computational resources to train.¹

1 INTRODUCTION

Self-attention-based architectures, in particular Transformers (Vaswani et al., 2017), have become the model of choice in natural language processing (NLP). The dominant approach is to pre-train on a large text corpus and then fine-tune on specific downstream tasks (Devlin et al., 2019). These models are trained on large datasets and are able to learn complex relationships between words of unprecedented size, with over 100B parameters (Brown et al., 2020; Lepikhin et al., 2020). With the model and training time scaling linearly with the number of words, there is no sign of diminishing returns.

In contrast, vision, however, convolutional architectures remain dominant (LeCun et al., 1989; Krizhevsky et al., 2012; He et al., 2016). Inspired by NLP successes, multiple works try combining CNN-like architectures with self-attention (Wang et al., 2018; Carion et al., 2020), some replacing the convolutional encoder (Ramanachandran et al., 2019; Wang et al., 2019). The latter models, while theoretically efficient, have not yet been able to effectively harness the power of self-attention due to the use of specialized attention patterns. Therefore, in large-scale image recognition, classic ResNet-like architectures are still state-of-the-art (Mahajan et al., 2018; Xie et al., 2020; Kolesnikov et al., 2019).

Inspired by the Transformer scaling successes in NLP, we experiment with applying a standard Transformer directly to images, with the fewest possible modifications. To do so, we split an image into patches and provide the sequence of linear embeddings of these patches as an input to a Transformer. Interestingly, this simple modification allows the model to learn the visual features (words) in an NLP application. We train the model on image classification in supervised fashion.

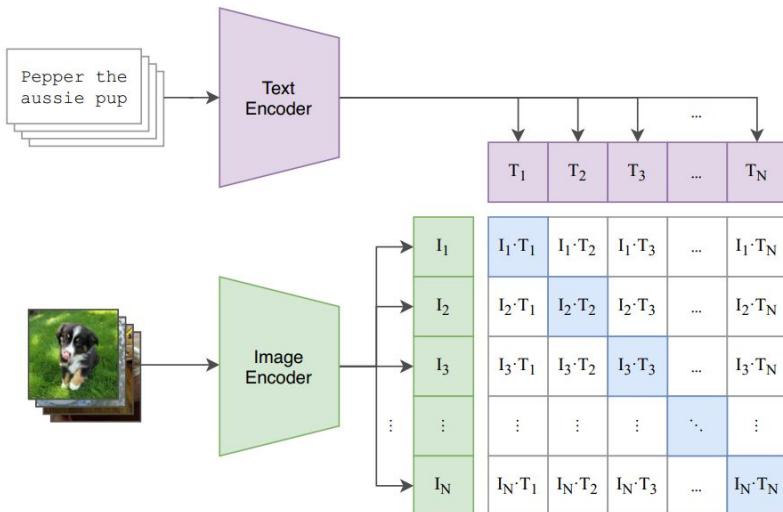
When trained on mid-sized datasets such as ImageNet without strong regularization, these models yield modest accuracies of a few percentage points below ResNets of comparable size. This seemingly discouraging outcome may be expected: Transformers lack some of the inductive biases

¹Fine-tuning code and pre-trained models are available at https://github.com/google-research/vision_transformer

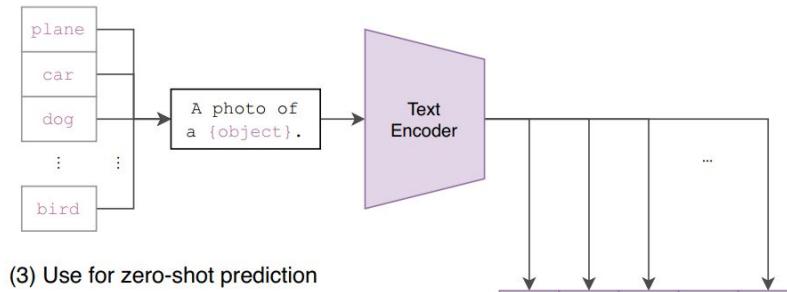
1

Dosovitskiy et al. (ICLR 2021)

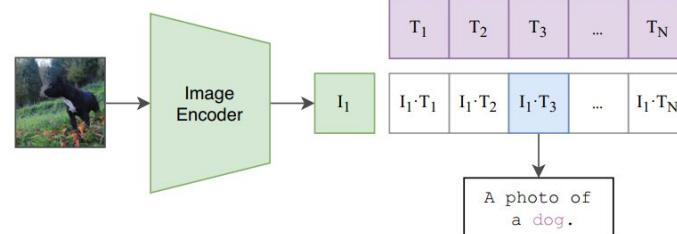
(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction





CLIP



Learning Transferable Visual Models From Natural Language Supervision

Alec Radford^{*1} Jong-Wook Kim^{*1} Chris Hallacy¹ Aditya Ramesh¹ Gabriel Goh¹ Sandhini Agarwal¹ Girish Sastry¹ Amanda Askell¹ Pamela Mishkin¹ Jack Clark¹ Gretchen Krueger¹ Dya Sutiskever¹

Abstract

State-of-the-art computer vision systems are trained to predict a fixed set of predetermined object categories. This restricted form of supervision limits their generality and utility since images are often required to be grouped by some other visual concept. Learning directly from raw text about images is a promising alternative which leverages a much broader source of supervision. We demonstrate that the standard supervised task of predicting which caption goes with which image is an efficient and scalable way to learn SOTA image captioning models. We train a dataset of 400 million (image, text) pairs collected from the internet. After pre-training, natural language is used to refine learned visual concepts (or “tokens”) and then to teach these models on downstream tasks. We study the performance of this approach by benchmarking on 30 different challenging computer vision datasets spanning tasks such as OCR, action recognition in videos, geo-localization, and many types of fine-grained object classification. The proposed approach is both fast and compact and is often competitive with a fully supervised baseline without the need for any dataset specific training. For instance, we match the accuracy of the state-of-the-art ResNet-50 model zero-shot without needing to use any of the 1.28 million training examples it was trained on. We release our code and pre-trained model weights at <https://github.com/openai/CLIP>.

1. Introduction and Motivating Work

Over 20 years ago Mori et al. (1999) explored improving content based image retrieval by training a model to predict the nouns and adjectives in text documents paired with images. Quantitative results demonstrated it was possible to learn more efficient image representations via manifold learning in the weight space of classifiers trained to predict words in captions associated with images. Srihari & Kakade (2004) later explored learning image representation learning by training multimodal Deep Boltzmann Machines on top of low-level image and text tag features.

Joulin et al. (2016) modernized this line of work by demonstrating that CNNs trained on image captioning tasks learn useful image representations. They converted the title, description, and hashing metadata of images in the ImageNet dataset (Thomee et al., 2010) into a bag-of-words multi-label classification problem and then they pre-trained AlexNet (Krizhevsky et al., 2012) to predict these learned representations which performed similarly well to a fully supervised pre-training on the same task. Radford et al. (2017) then extended this approach to predicting phrase n-grams in addition to individual words and demonstrated the ability of their system to zero-shot transfer to other image

arXiv:2103.00020v1 [cs.CV] 26 Feb 2021

arXiv:2309.16671v5 [cs.CV] 28 Dec 2024

Radford et al. (ICML 2021)

Preprint

DEMYSTIFYING CLIP DATA

Hu Xu¹ Saining Xie^c Xiaojing Ellen Tan^b Po-Yao Huang^b Russell Howes^b Vasu Sharma^b
Shang-Wen Li^d Gergi Ghosh^e Luke Zettlemoyer^{f,g} Christoph Feichtenhofer^c
^{FAIR, Meta AI} ^{New York University} ^{University of Washington}

ABSTRACT

Contrastive Language-Image Pre-training (CLIP) is an approach that has advanced research and applications in computer vision, fueling state-of-the-art systems for image captioning, image retrieval, and image generation. A key aspect of CLIP is that it is *data and not the model* architecture or pre-training objective. However, CLIP only provides very limited information about its data and how it has been curated. In this work, we reveal the details of CLIP’s data curation and its training with its model parameters. In this work, we intend to reveal CLIP’s data curation approach and in our pursuit of making it open to the community introduce new metrics for evaluating the quality of datasets. Our work shows that CLIP takes a raw public and meta-data (derived from CLIP’s concepts) and yields a balanced subset over the metadata distribution. Our experimental study rigorously isolates the effect of the data size and quality on CLIP’s performance. When we apply CLIP to CommonCrawl with 400M image-text data pairs outperforms CLIP’s data on multiple standard benchmarks. In zero-shot ImageNet classification, MetaCLIP with 100M images and 100M captions, trained with a budget of 100 hours totaling to 1B data, while maintaining the same training budget, attains 72.4%. Our observations hold across various model sizes, exemplified by ViT-b/16 producing 82.1%. Our full code and training data distribution over metadata is available at <https://github.com/lemoncode/research/MEAC-CLIP>.

1 INTRODUCTION

Deep learning has revolutionized the field of artificial intelligence, and pre-trained models have played a pivotal role in democratizing access to cutting-edge AI capabilities. However, the training data used to create these models is often concealed from the public eye, shrouded in secrecy.

The increasing availability of pre-trained models for public use contrasts sharply with the lack of transparency regarding their training data. Further, proprietary concerns, such as copyright issues, often prevent researchers from accessing training data, thus hindering efforts to explore new approaches for curating high-quality training data that can be shared openly.

In the vision-language domain, the dominant model and learning approach is Contrastive Language-Image Pre-training (CLIP) (Radford et al., 2021), a simple technique to learn from image-text pairs. CLIP uses a large dataset of image-text pairs from the web, specifically the WIT400M dataset which is curated from the web. Despite its popularity, the specifics of CLIP’s curation process have remained a mystery, captivating the research community for years.

Follow-up works (Schuhmann et al., 2022; 2021) have attempted to replicate CLIP’s data, but with a modified curation strategy. While these approaches use a different unknown data source and curation methodology, these approaches remove noise by applying the CLIP model as a hard blackbox filter which in turn is a form of distilling WIT400M information captured in CLIP.

The advantages of CLIP’s curation are apparent. First, it starts from scratch, bypassing the introduction of noise from a pre-existing dataset. Second, it preserves the original data structure after metadata, maximizing signal preservation while mitigating, rather than removing, noise in the data.¹ Such distribution lays the groundwork for task-agnostic data, a crucial part of foundation models.

¹For example, a filter on digits can remove noise from date or id strings but remove signal for tasks that involve OCR (e.g., MNIST), or a filter removing text with less than 5 characters can remove signal “dog”.

Hu Xu et al. (ICLR 2024)

Transfer Learning

```
1  #
2  #
3  #
4  import transformers
5
6
7  #
8  #
9  #
10 c = transformers.CLIPModel.from_pretrained('openai/clip-vit-base-patch32')
11 v = transformers.ViTModel.from_pretrained('google/vit-base-patch16-224-in21k')
12 params = lambda m: sum(p.numel() for p in m.parameters())
13
14
15 #
16 #
17 #
18 print("CLIP:", params(c)) # 151,277,313
19 print("ViT:", params(v)) # 86,389,248
20
21
22 #
23 #
24 #
25 print("CLIP:", c)
26 print("ViT:", v)
27
```

```
ViTModel(  
    (embeddings): ViTEmbeddings(  
        (patch_embeddings): ViTPatchEmbeddings(  
            (projection): Conv2d(3, 768, kernel_size=(16, 16), stride=(16, 16))  
        )  
        (dropout): Dropout(p=0.0, inplace=False)  
    )  
    (encoder): ViTEncoder(  
        (layer): ModuleList(  
            (0-11): 12 x ViTLayer(  
                (attention): ViTSdpAttention(  
                    (attention): ViTSdpSelfAttention(  
                        (query): Linear(in_features=768, out_features=768, bias=True)  
                        (key): Linear(in_features=768, out_features=768, bias=True)  
                        (value): Linear(in_features=768, out_features=768, bias=True)  
                        (dropout): Dropout(p=0.0, inplace=False)  
                    )  
                    (output): ViTSdpSelfOutput(  
                        (dense): Linear(in_features=768, out_features=768, bias=True)  
                        (dropout): Dropout(p=0.0, inplace=False)  
                    )  
                )  
                (intermediate): ViTIntermediate(  
                    (dense): Linear(in_features=768, out_features=3072, bias=True)  
                    (intermediate_act_fn): GELUActivation()  
                )  
                (output): ViTOuput(  
                    (dense): Linear(in_features=3072, out_features=768, bias=True)  
                    (dropout): Dropout(p=0.0, inplace=False)  
                )  
                (layernorm_before): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
                (layernorm_after): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
            )  
        )  
        (layernorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
    )  
    (pooler): ViTPooler(  
        (dense): Linear(in_features=768, out_features=768, bias=True)  
        (activation): Tanh()  
    )  
)
```



FOUNDERS AND CODERS

CLIP

```
CLIPModel(
    (text_model): CLIPTextTransformer(
        (embeddings): CLIPTextEmbeddings(
            (token_embedding): Embedding(49408, 512)
            (position_embedding): Embedding(77, 512)
        )
        (encoder): CLIPEncoder(
            (layers): ModuleList(
                (0-11): 12 x CLIPEncoderLayer(
                    (self_attn): CLIPSdpAttention(
                        (k_proj): Linear(in_features=512, out_features=512, bias=True)
                        (v_proj): Linear(in_features=512, out_features=512, bias=True)
                        (q_proj): Linear(in_features=512, out_features=512, bias=True)
                        (out_proj): Linear(in_features=512, out_features=512, bias=True)
                    )
                    (layer_norm1): LayerNorm((512,), eps=1e-05, elementwise_affine=True)
                    (mlp): CLIPMLP(
                        (activation_fn): QuickGELUActivation()
                        (fc1): Linear(in_features=512, out_features=2048, bias=True)
                        (fc2): Linear(in_features=2048, out_features=512, bias=True)
                    )
                    (layer_norm2): LayerNorm((512,), eps=1e-05, elementwise_affine=True)
                )
            )
        )
        (final_layer_norm): LayerNorm((512,), eps=1e-05, elementwise_affine=True)
    )
    (vision_model): CLIPVisionTransformer(
        (embeddings): CLIPVisionEmbeddings(
            (patch_embedding): Conv2d(3, 768, kernel_size=(32, 32), stride=(32, 32), bias=False)
            (position_embedding): Embedding(50, 768)
        )
        (pre_layernorm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
        (encoder): CLIPEncoder(
            (layers): ModuleList(
                (0-11): 12 x CLIPEncoderLayer(
                    (self_attn): CLIPSdpAttention(
                        (k_proj): Linear(in_features=768, out_features=768, bias=True)
                        (v_proj): Linear(in_features=768, out_features=768, bias=True)
                        (q_proj): Linear(in_features=768, out_features=768, bias=True)
                        (out_proj): Linear(in_features=768, out_features=768, bias=True)
                    )
                    (layer_norm1): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
                    (mlp): CLIPMLP(
                        (activation_fn): QuickGELUActivation()
                        (fc1): Linear(in_features=768, out_features=3072, bias=True)
                        (fc2): Linear(in_features=3072, out_features=768, bias=True)
                    )
                    (layer_norm2): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
                )
            )
        )
        (post_layernorm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
    )
    (visual_projection): Linear(in_features=768, out_features=512, bias=False)
    (text_projection): Linear(in_features=512, out_features=512, bias=False)
)
```



Machine
Learning
Institute

Datasets

Dataset Viewer

Auto-converted to Parquet | API | View in Dataset Viewer

Split (1)
test · 31k rows

Search this dataset

image	caption	sentids	split	img_id	filename
image	list	list	string · classes	string · lengths	string · lengths
[image]	["Two young guys with shaggy hair look at..."	["0", "1", "2", "3", -]	train	0	1000092795.jpg
[image]	["Several men in hard hats are operating a..."	["5", "6", "7", "8", -]	train	1	10002456.jpg
[image]	["A child in a pink dress is climbing up..."	["10", "11", "12", -]	train	2	1000268201.jpg
[image]	["Someone in a blue shirt and hat is..."	["15", "16", "17", -]	train	3	1000344755.jpg
[image]	["Two men, one in a ..."	["20", -]			

< Previous 1 2 3 ... 311 Next >

Flickr30k

Making the V in VQA Matter:
Elevating the Role of Image Understanding in Visual Question Answering

Yash Goyal¹ Teja Khot¹ Douglas Summers-Stay² Dhruv Batra³ Dev Parikh³
¹Virginia Tech, ²Facebook AI Research, ³Georgia Institute of Technology
arXiv:1612.00837v3 [cs.CV] 15 May 2017

Abstract

Problems in the intersection of vision and language are of significant interest both for their inherent difficulty and for the rich set of applications they enable. In this paper, we propose a new dataset, SlideVQA, where language tends to be a simpler signal for learning than visual modalities, resulting in models that ignore visual information and fail on simple questions.

We propose to counter these language priors for the task of Visual Question Answering (VQA) by creating a balanced VQA dataset. Specifically, we balance the popular VQA dataset [1] by collecting complementary images such that every image has two different answers, not just one, and not just a single image, but rather a pair of similar images that results in two different answers to the question. Our proposed dataset, SlideVQA, contains 2.8M images and 2.8M questions, and is publicly available at <http://slidevqa.csail.mit.edu> as part of the 2nd iteration of the Visual Question Answering Benchmark.

We further benchmark a number of state-of-art VQA models on our balanced dataset. All models perform significantly better on our dataset than on the original VQA dataset. These models have indeed learned to exploit language priors, and this is reflected in the fact that they consistently outperform the VQA baseline on our dataset.

Our data collection protocol for identifying complementary images enables us to develop a novel interpretable model, which in addition to providing an answer to the question, also provides a detailed, step-by-step example based explanation. Specifically, it identifies an image that is similar to the original image, but it believes has a different answer. This can help users “trust” the answer to the question “Is there a person in the picture?” on images actually containing clock towers. As is particularly perverse example – for questions

The first two authors contributed equally.

Figure 1: Examples from our balanced VQA dataset.

SlideVQA: A Dataset for Document Visual Question Answering on Multiple Images

Ryota Tanaka, Kyosuke Nishida, Koukei Nishida, Taku Hasogawa, Isamu Saito, Kuniko Saito
NTT Human Information Laboratory, NTT Corporation
{ryota.tanaka.ng, kyosuke.nishida.e, koukei.nishida.e, taku.hasogawa.e, isamu.saito.d, kuniko.saito.k}@hilo.att.co.jp

Visual question answering on document images that contain textual, visual, and layout information, called document visual question answering (DVQA), have been proposed for developing document VQA systems. DVQA datasets have been proposed for developing document VQA systems. However, most of them have been proposed for document images with a single image and not for multiple images. We propose a new dataset, SlideVQA, containing 2.8M images and 2.8M questions, for DVQA. The dataset consists of a slide deck composed of multiple slide images and a corresponding set of questions. Each slide deck contains a set of images and answers to the question. Slide decks are one of the most common document formats used in business and educational environments for communication. As shown in Figure 1, SlideVQA requires complex reasoning over slide images, including reading and understanding the visual content. These reasoning skills play essential roles in MRC tasks (Yang et al., 2016).

Our main contributions are summarized as follows:

- We introduce a novel task and dataset, SlideVQA, which requires complex reasoning over slide decks and comprehend a slide deck. It is the largest multi-image document VQA dataset consisting of 2.8M images and 2.8M questions. The dataset is composed of 20 slide decks and 14.3k questions. It also provides bounding boxes around textual and visual elements in each slide deck, which are useful for arithmetic expressions for numerical reasoning.
- We develop a Multi-Modal Multi-image Document VQA (MMDVQA) system for document image selection and question answering tasks and to enhance numerical reasoning by generating arithmetic expressions.

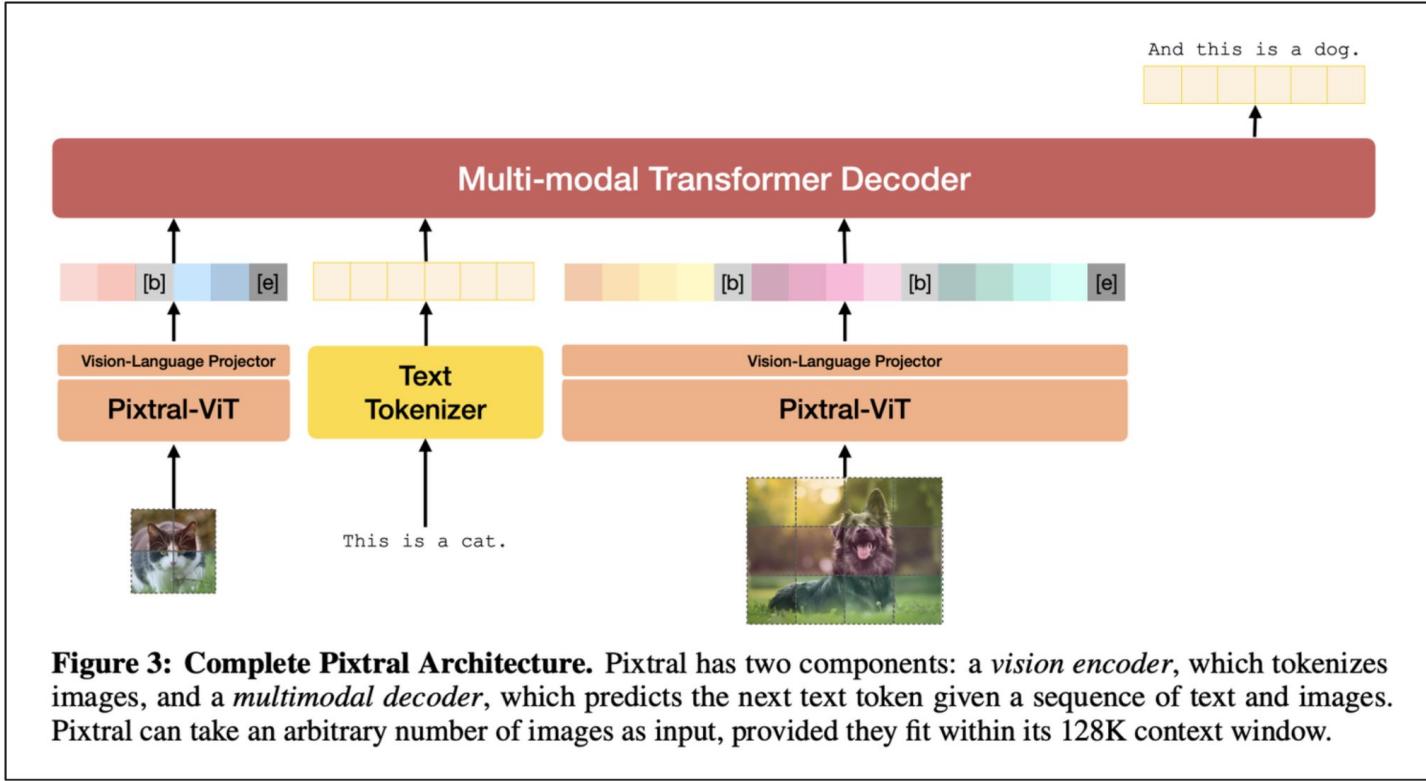
*Our dataset and codes are publicly available at <https://github.com/ntt-mmdu/slidevqa>.

SlideVQA

arXiv:2301.04883v1 [cs.CL] 12 Jan 2023

VQAv2

Pixtral 12B



Task

- use the hidden state from ViT or CLIP
- code only the decoder
- use notebooks for quick experiments
- deploy to Hetzner

Extra

- use “EleutherAI/pythia-160m” pre-trained



Good luck!

Teams

Recurrent Rebels

Aygun
Stanley
James
Josh

Gradient Gigglers

Dimitar
Andrea
Evelyn
Emil

Overfitting Overlords

Liam
David
Guillaume
Ollie

Perceptron Party

Loredana
Dimitris
Pyry
Amy

Backprop Bunch

Kori
Ardrit
Neville
Gaurav

Dropout Disco

Kenton
Yurii
Artemis
Nnamdi

Activation Aces

Milo
Daniel
Filippo
Coline