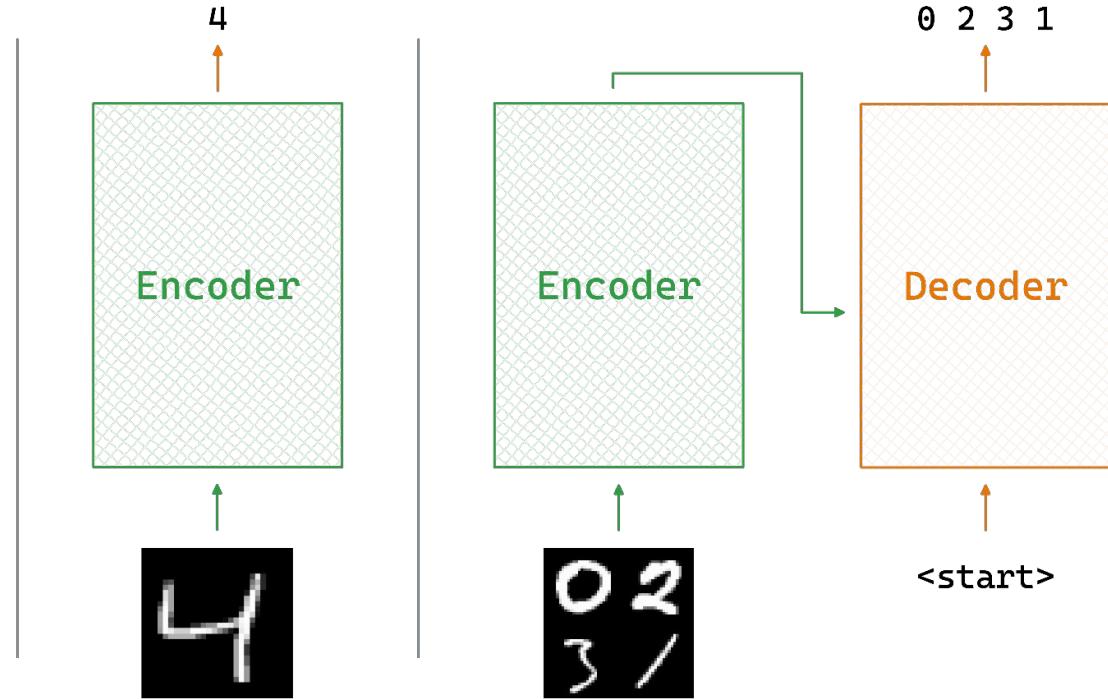


# Multimodal transfer learning

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

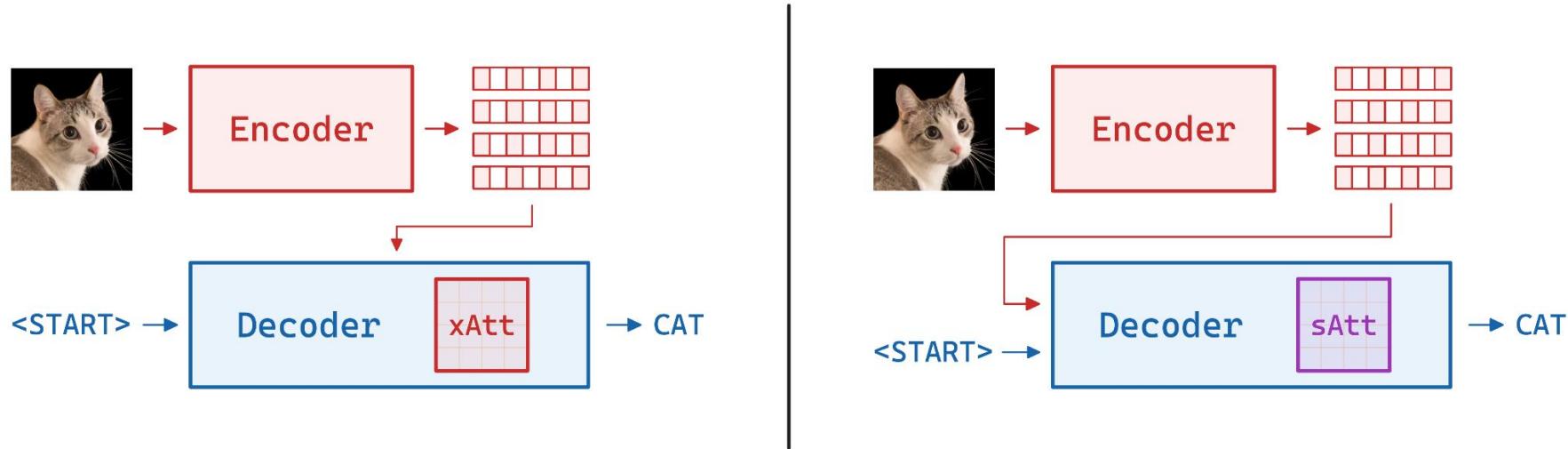


# Task

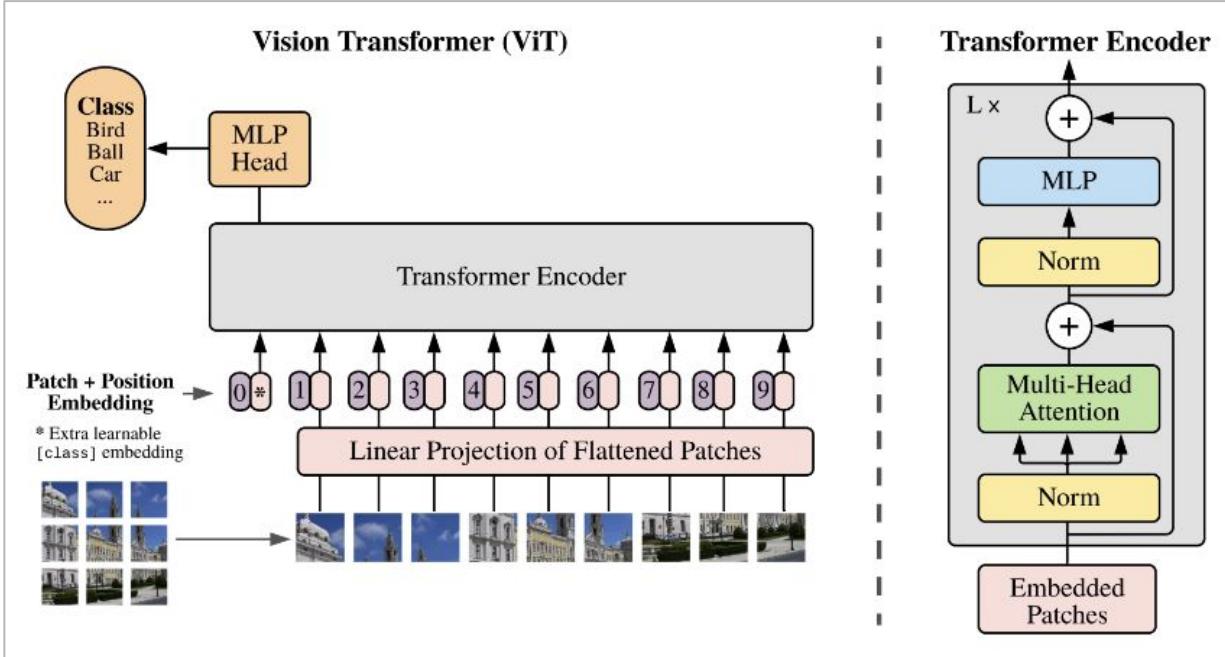


A little girl  
climbing into a  
wooden playhouse.

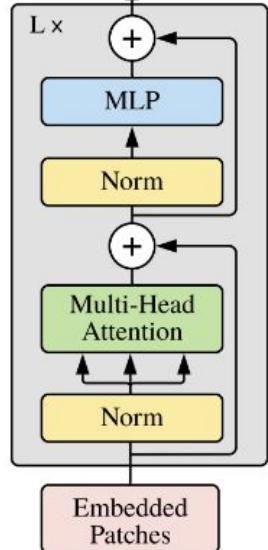
# Transfer Learning



# 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 Dosovitskiy<sup>1</sup>, Lucas Beyer<sup>2</sup>, Alexander Kolesnikov<sup>3</sup>, Dirk Weissenborn<sup>2</sup>, Xiaohua Zhai<sup>2</sup>, Thomas Unterthiner<sup>2</sup>, Matthias Minderer<sup>2</sup>, Matthias Minderer<sup>2</sup>, Georg Heigold<sup>2</sup>, Sylvain Gelly<sup>2</sup>, Jakob Uszkoreit<sup>2</sup>, Neil Houlsby<sup>2</sup>  
<sup>1</sup>Equal technical contribution  
<sup>2</sup>Google 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 achieves state-of-the-art results on ImageNet, CIFAR-100, VTFB, etc. Our Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring substantially fewer computational resources to train.<sup>1</sup>

## 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 scale growing, so does the cost of training.

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 (Ramanathan 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.

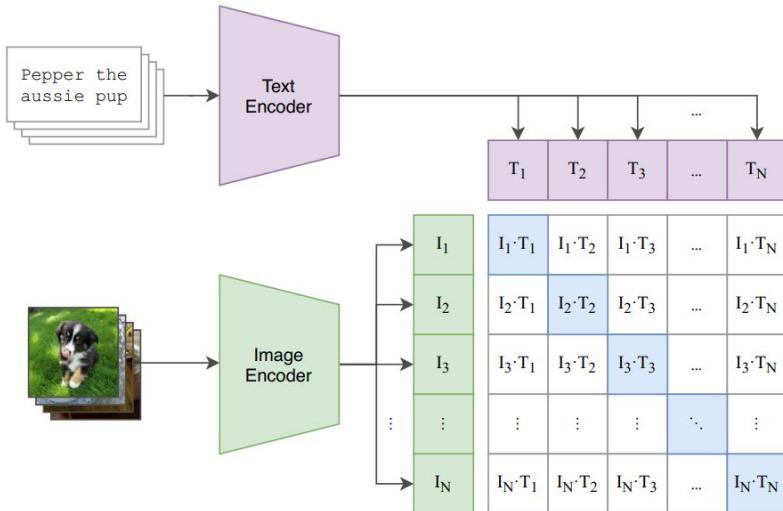
When trained on mid-sized datasets such as ImageNet without strong regularization, these models yield modest accuracies of a few percentage points below ResNets of comparable size. This seemingly discouraging outcome may be expected: Transformers lack some of the inductive biases

<sup>1</sup>Fine-tuning code and pre-trained models are available at [https://github.com/google-research/vision\\_transformer](https://github.com/google-research/vision_transformer)

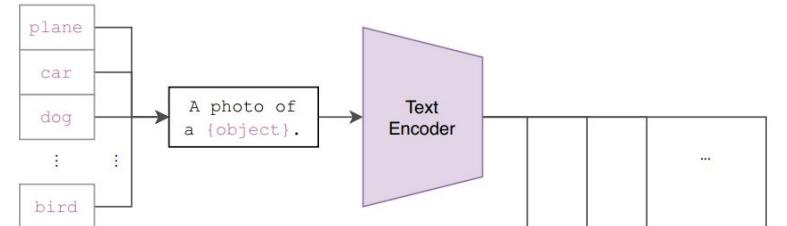
1

Dosovitskiy et al. (ICLR 2021)

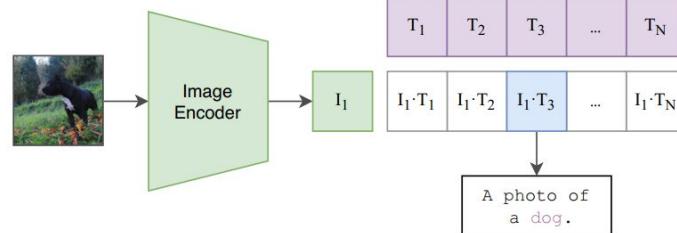
(1) Contrastive pre-training



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction





# CLIP



## Learning Transferable Visual Models From Natural Language Supervision

Alec Radford<sup>\*1</sup> Jong-Wook Kim<sup>\*1</sup> Chris Hallacy<sup>1</sup> Aditya Ramesh<sup>1</sup> Gabriel Goh<sup>1</sup> Sandhini Agarwal<sup>1</sup> Girish Sastry<sup>1</sup> Amanda Askell<sup>1</sup> Pamela Mishkin<sup>1</sup> Jack Clark<sup>1</sup> Gretchen Krueger<sup>1</sup> Dya Sutiskever<sup>1</sup>

### 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 labeled to predict other visual concepts. 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”) in the model. This allows us to transfer to downstream datasets involving the specific visual concepts or dataset specific curations. Pretrained systems like GPT-3 (Brown et al., 2020) are now competitive across many tasks with bespoke models while requiring little to no dataset engineering.

These results suggest that the aggregate supervision accessible to modern pre-training methods within web-scale collections of text surpasses that of high-quality crowd-labeled NLP datasets. However, in other domains such as vision, it is not clear how to transfer learned models on crowd-labeled datasets such as ImageNet (Deng et al., 2009). Could scalable pre-training methods which learn directly from raw text offer a similar breakthrough in computer vision? Prior work is encouraging.

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) proposed a similar approach to representation learning by training multimodal Deep Boltzmann Machines on top of low-level image and text tag features. Jodoin et al. (2010) formalized this line of work and demonstrated that CNNs trained on image-caption pairs learn useful image representations. They converted the title, description, and hashing metadata of images in the CIFAR100 dataset (Krizhevsky et al., 2009) 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 baseline 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

### 1 Introduction and Motivating Work

Pre-training methods which learn directly from raw text have revolutionized NLP over the last few years (Dai & Le, 2015; Peters et al., 2018; Howard & Ruder, 2018; Radford et al., 2018; Devlin et al., 2018; Raffel et al., 2019).

<sup>\*</sup>Equal contribution. <sup>1</sup>OpenAI, San Francisco, CA 94110, USA  
Correspondence to: {alec, jongwook}@openai.com>

Radford et al. (ICML 2021)

Preprint

## DEMYSTIFYING CLIP DATA

Hu Xu<sup>1</sup> Saining Xie<sup>2</sup> Xiaojing Ellen Tan<sup>1</sup> Po-Yao Huang<sup>1</sup> Russell Howes<sup>1</sup> Vasu Sharma<sup>1</sup> Shang-Wen Li<sup>1</sup> Gergi Ghosh<sup>1</sup> Luke Zettlemoyer<sup>1,3</sup> Christoph Feichtenhofer<sup>1</sup>  
FAIR, Meta AI <sup>1</sup>New York University <sup>2</sup>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. The key insight of CLIP is that *it is the 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 curation process of CLIP’s data 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 benchmarks. Our experiments show that CLIP’s curation approach makes it a raw data pool 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 curation process and shows that when the curation is applied to CommonCrawl with 400M image-text data pairs outperforms CLIP’s data on multiple standard benchmarks. In zero-shot ImageNet classification, MetaCLIP (our curation approach) achieves 74.4% top-1 accuracy on ImageNet-1k with only 100M training images, while maintaining the same training budget, attains 72.4%. Our observations hold across various model sizes, exemplified by ViT-16@G producing 82.1% top-1. Our full code and training data distribution over metadata is available at <https://github.com/lemoncode/research/MetaCLIP>.

### 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 the release of training data, which motivates data scientists to explore new approaches for curating high-quality training data that can be shared openly arises.

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 model-free approach to curating data. While these approaches remove noise from the data and 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 the training data. Second, it preserves the original data structure, including the meta-data, maximizing signal preservation while mitigating, rather than removing, noise in the data. Such distribution lays the groundwork for task-agnostic data, a crucial part of foundation models.

<sup>1</sup>For 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”.

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

Hu Xu et al. (ICLR 2024)

1

# 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
```



# ViT



```
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): ViTSdpaAttention(  
                    (attention): ViTSdpaSelfAttention(  
                        (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): ViTSelfOutput(  
                        (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()  
    )  
)
```



# CLIP



```
CLIPModel(  
    (text_model): CLIPTextTransformer(  
        (embeddings): CLIPTextEmbeddings(  
            (token_embedding): Embedding(49408, 512)  
            (position_embedding): Embedding(77, 512)  
        )  
        (encoder): CLIPTextEncoder(  
            (layers): ModuleList(  
                (b-11): 12 x CLIPTextEncoderLayer(  
                    (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): CLIPVisionEncoder(  
            (layers): ModuleList(  
                (b-11): 12 x CLIPVisionEncoderLayer(  
                    (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)  
    )
```

# 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
[img]	[ "Two young guys with shaggy hair look at..."]	[ "0", "1", "2", "3", ...]	train	0	1000092795.jpg
[img]	[ "Several men in hard hats are operating a..."]	[ "5", "6", "7", "8", ...]	train	1	10002456.jpg
[img]	[ "A child in a pink dress is climbing up ..."]	[ "10", "11", "12", ...]	train	2	1000268201.jpg
[img]	[ "Someone in a blue shirt and hat is..."]	[ "15", "16", "17", ...]	train	3	1000344755.jpg
[img]	[ "Two men, one in a ..."]	[ "20", ...]			

< Previous 1 2 3 ... 311 Next >

arXiv:1612.00837v3 [cs.CV] 15 May 2017

Flickr30k

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

Yash Goyal<sup>1</sup> Teja Khot<sup>1</sup> Douglas Summers-Stay<sup>2</sup> Dhruv Batra<sup>3</sup> Dev Parikh<sup>3</sup>  
<sup>1</sup>Virginia Tech, <sup>2</sup>AMP Research Laboratory, <sup>3</sup>Georgia Institute of Technology  
{yogoyal, tjkhot}@vt.edu, {douglas.s.summers-stay, cbatra, parikh}@gatech.edu

**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 argue that the visual modality in natural language tends to be a simpler signal for learning than visual modalities, resulting in models that ignore visual information and fail to learn from it.

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, called *SlideVQA*, contains 2.8M questions and 1.4M images, 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. This demonstrates the first concrete empirical evidence of what seems to be a well-known phenomenon in machine learning.

Finally, 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. For example, “What is the tower in the picture?” on images actually containing clock towers. As is particularly perverse example – for questions

“Who is waving glasses?”, “Is the umbrella upside down?”, “How many children are in the bed?”

Figure 1: Examples from our balanced VQA dataset.

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

VQAv2

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

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

In reason about document layout, textual content, and visual elements (Mathew, Karatzas, and Jawahar 2021; Tanaka, Nishida, and Saito 2022). The primary task is to extract the primary context in a document to text (e.g., e-mails and forms) and the task is to understand it on the basis of the context. SlideVQA is the first dataset that has achieved nearly human-level performance (Xu et al. 2021; Dong et al. 2022). However, the performance is still far from human level when it comes to handling diverse real-world documents. One of the reasons is that current models are not capable of performing reasoning across images. Most of the existing datasets focus on testing reasoning ability on a single image (Yang et al. 2019; Xu et al. 2021). VQA and VQA2 models still have trouble understanding documents that contain multiple images. This is because they do not require numerical reasoning (Mathew et al. 2022).

To address the above challenges, we introduce a new document visual question answering (DVQA) task. DVQA given a slide deck composed of multiple slide images and a corresponding question, asks the model to find the most relevant images and answers to the question. Slide decks are one of the most common ways to communicate ideas and visual elements for communication. As shown in Figure 1, SlideVQA requires complex reasoning over slide images, including reasoning over multiple images. Therefore, reasoning skills play essential roles in MRC tasks (Yang et al. 2022).

Our main contributions are summarized as follows:

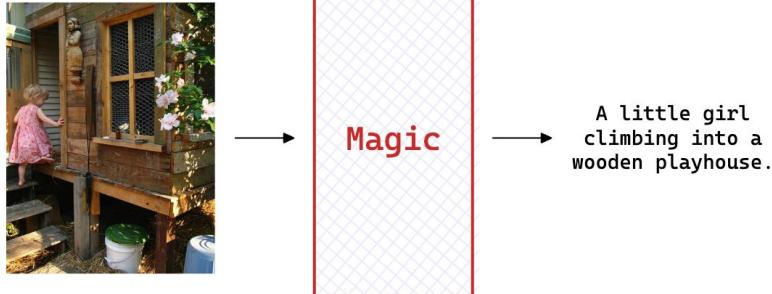
- We introduce a novel task and dataset, SlideVQA, which requires reasoning over multiple images and comprehend a slide deck. It is the largest multi-image document VQA dataset, consisting of 20 million images and 14.5M questions. It also provides bounding boxes around textual and visual elements in each image, which is useful for arithmetic expressions for numerical reasoning.
- We develop a Multi-Modal Multi-image Document Visual Question Answering (DVQA) task for document selection and question answering tasks to enhance numerical reasoning by generating arithmetic expressions.

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

SlideVQA

# Task

1. learn how to use the hidden state from ViT or CLIP
2. code from scratch the decoder
3. create synthetic datasets and save them on Hugging Face
  - a. **get used to it!**
  - b. use Qwen2.5-VL-3B-Instruct
4. train and experiment with datasets and alignment
5. **use “Qwen/Qwen3-0.6B-Base” pre-trained as decoder**



# Good luck!

## Recurrent Rebels

Ben Lioong  
Kadriye Turkcan  
Rosh Beed  
Joao Esteves

## Gradient Gigglers

Jingyan Chen  
Umut Sagir  
Tyrone Nicholas  
James Yan

## Overfitting Overlords

Tao Zamorano  
Andrew  
Melanie Wong  
Felipe Lavratti

## Hyperparameter Hippies

Miguel Parracho  
Prima Gouse  
Esperanza Shi  
Rasched Haidari

## Perceptron Party

Yali Pan  
Andrei Zhirnov  
Dan Goss  
Adam Beedell

## Backprop Bunch

Ben Bethell  
Arjuna James  
Helen Zhou  
Maria Sharif

## Dropout Disco

Hikaru Tsujimura  
Ewan Beattie  
Ethan Edwards  
Nikolas Kuhn

## Kernel Kittens

Jacob Jenner  
Marcin Tolysz  
Tomas Krajcoviech  
Aparna Pillai

## Bayesian Buccaneers

Anton Dergunov  
David Edev  
James Carter  
Peter O'Keefe

## Feature Fiestas

Ben Williams  
Charles Cai  
Clement Ha