# Prompted Multi-Modal learning for Vision-Language Understanding for Image Caption Generation

**Mukesh Bangalore Renuka**
Columbia University
`mb4862@columbia.edu`

**Yogesh Patodia**
Columbia University
`yp2607@columbia.edu`

## Abstract

Image captioning systems have been of significant interest in the Machine Learning community since the evolution of Deep learning, specially in the field of Vision-Language understanding. The aim of these systems are to predict an informative caption about an input image. However traditional deep learning systems are resource intensive, lack variations in these generations and also lacks the ability to be prompted as to the style of captioning. In this paper we propose a novel training mechanism to address this issue of prompting using dual vision text encoders. We also demonstrate it's ability in solving Visual Question Answering and other related tasks in the field of Multimodal Vision-Language learning.

## 1 Introduction

Computer Vision (CV) and Natural Language Processing(NLP) are two fundamental parts of Artificial Intelligence which focus on understanding image and language respectively. From the advent of these two factions of AI, researchers have found ways to marry these two schools in order to build more intelligent systems, thus forming one of the forms of Multi Modal learning.

Image Captioning represents the fundamental task in this field. The task entails providing a meaningful caption in natural language which describes the input image. Although it seems rather trivial for humans, it poses few challenges for AI systems. The main ones being semantic understanding and variability of captions. Semantic Understanding involves detecting objects in the image and understanding the different relations between the objects in the image. Variability of captions deals with the different ways people can meaningfully interpret and caption an image. Although a picture is worth a thousand words (or 16*16 words if you
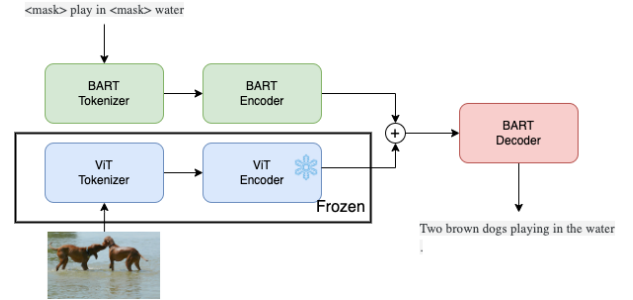


Figure 1: Proposed Image captioning pipeline with Prompting

subscribe to the Visual Transformers paper), these words might be different for different people.

In our project, we aim to leverage transformer based models such as BART and ViT model pretrained on very large datasets to perform prompted image captioning. We wish to give the model the ability to be prompted as this can enable the model to perform a variety of tasks apart from Image captioning such as Visual Question Answering, Visual Entailment, Visual Reasoning etc.

## 2 Related Works

Image captioning has been extensively studied in literature. Most deep learning approaches (Stefanini et al., 2021) tend to use visual features generated by traditional Deep Learning CNN based systems and textual decoders to produce the desired output. Early approaches (Chen and Zitnick, 2014; Chen et al., 2016; Fang et al., 2014) relied on visual embeddings generated by traditional CNNs on various tasks whereas latter approaches (Anderson et al., 2017; Li et al., 2020; Zhou et al., 2020) relied on more sophisticated object detection models to generate visual features. However the decoder must bridge the gap between textual representations and Visual representation in order to generate good captions. This involves various

LSTM based models (Chen et al., 2018; Vinyals et al., 2014) to decode these features and generate these captions in Natural Language.

This led to increase in computational requirements. ClipCap (Mokady et al., 2021) addresses this issue by using pretrained deep learning CLIP model (Contrastive Language-Image Pre-Training) as the encoder, to simplify the captioning pipeline while improving the performance. However this still lacks the ability to prompt networks

The BART paper (Lewis et al., 2019) which pretrains a model combining Bidirectional and Auto-Regressive Transformers and a denoising autoencoder in order to pretrain large scale sequence to sequence models. In their paper, they mention several masking schemes such as token masking, token deletion, text infilling, sentence permutations and document rotation in order to obtain a more robust pretraining regiment.

## 3 Data

We train our model **Flickr30K** dataset (Plummer et al., 2015) which contains around 30K images. Each image contains 5 captions. Thus, we have around **150K** samples. The dataset contains images and 5 captions for each image describing it.

We plan on using The Microsoft Common Objects in Context (**COCO**) captions dataset (Chen et al., 2015). Coco is a large-scale object detection, segmentation, and captioning dataset. We primarily focus on the image captioning part of the dataset which contains 413,915 captions for 82,783 images in training, 202,520 captions for 40,504 images in validation and 379,249 captions for 40,775 images in test set. We will use this for experiments in the latter half of the semester.

We also plan on evaluating the VQA performance of this model using the VQA 2.0 dataset (Goyal et al., 2017) which is a large dataest of images, related questions and their answers. It consists of 82,783 images, having 443,757 questions and 4,437,570 answers as the training set. We hope to subsample this dataset to ensure we can better test the model with the compute constraints.

## 4 Methodology

### 4.1 Model Architecture

The overview of the model architecture is shown in fig 2. We make use of the Visual Transformer (Dosovitskiy et al., 2020) to extract powerful representations of the input image, which is then con-
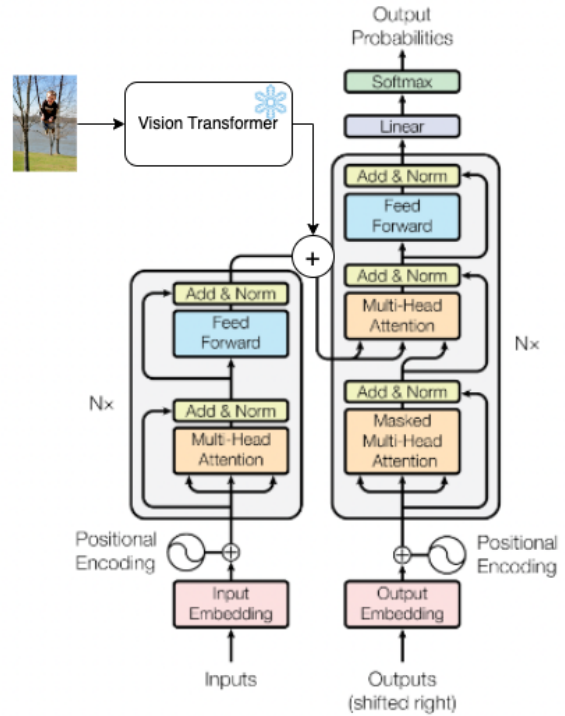


Figure 2: Architecture Diagram, illustration borrowed from (Vaswani et al., 2017)

catenated to the encoder representations from the BART encoder (Lewis et al., 2019) before being fed to the BART decoder (Lewis et al., 2019) . The attention mask corresponding to the image embeddings were set to be 1.

**ViT Image Encoder** Vision Transformer or ViT (Dosovitskiy et al., 2020) is a transformer only Image understanding model which can be used to perform various tasks such as image classification, segmentation and various other tasks in Computer Vision. It involves breaking up the input image into 16x16 patches which are then directly fed into a transformer based model along with positional embeddings. The ViT model was pretrained on ImageNet-21k, a dataset consisting of 14 million images and 21k classes, and fine-tuned on ImageNet, a dataset consisting of 1 million images and 1k classes.

**BART** BART (Lewis et al., 2019) is a model combining Bidirectional and Auto-Regressive Transformers and a denoising autoencoder in order to pretrain large scale sequence to sequence models. Several masking schemes such as token masking, token deletion, text infilling, sentence permutations and Document rotation are applied in order to obtain a more robust pretraining regiment. BART can be seen as a generalization of BERT (with its bidi-

| Component | Input | Output |
|---|---|---|
| Vision Encoder (ViT) | \<Input Image\> | \<Image Encoding\> |
| BART Text Encoder | \<Masked Caption\> | \<Text Encoding\> |
| BART Joint Decoder | \<Image Encoding\>\<Text Encoding\> | \<Output IDs\> |

Table 1: Descriptions of Inputs and outputs to our model.

rectional encoder) and GPT (with its left to right decoder).

## 4.2 Masking Strategies

We provide the masked caption as the input to the model and hope to wean of the dependency of the model on the provided input text. However we hope to teach BART to retain the ability to rely on the input text if/when it is present so that the model can be prompted to generate a suitable caption. We broadly use three different masking strategy which are directly inspired by BART:

- **Masking the entire caption**: Where we mask the entire caption and provide the model with just an empty string as the input.

- **Token Masking**: Where we mask a subset of tokens with a placeholder "\<mask\>". We implement two versions of such a masking scheme, namely random token masking and epoch aware token masking.

- **Text Infilling**: Where we replace chunks/phrases of the caption with a single mask token before feeding it as input. From the BART paper, we see the utility in such a type of masking.

**Describing the Inputs** The inputs and outputs of all components during training are described in table 1

## 5 Experiments

In this section, we share experimental and implementation details regarding the pretraining and pre-finetuning steps of our approach.

## 5.1 Experimental Design

We use all models from huggingface. The Vision Encoder model we use was the the ViT model that was pretrained on ImageNet-21k and fine-tuned on the larger ImageNet. The ViT model consisted of 86.39M parameters. These weights were frozen and not finetuned. The BART model was the base model which had 12 layers each in encoder and decoder layer, 1024-hidden and 140M parameters.

The model was pretrained of the Masked Language Modelling task on massive corpora of text such as the BOOKCORPUS dataset, WIKIPEDIA dataset, CC-NEWS dataset and the OPENWEB-TEXT dataset. The BART base model consisted of 139.42M parameters.

We trained the models using the AdamW optimizer with an initial learning rate of 0.001 and with a exponential Leearning Rate scheduler with a gamma of 0.9 on a singular NVIDIA A100 GPU.

## 5.2 Training Stratergies

We froze the image encoder to ensure reduced computational overhead. We also added early stopping and checkpointing to make sure the model does not over fit and more importantly to conserve compute resources. We then established three masking methodologies whose specifics we share here.

**Baseline (Entire caption masking)** To establish a baseline to measure the improvements in performance due to pretraining strategies, we feed an empty input string as the prompt to the caption. We hope that it performs akin to a vanilla Image captioning system. While training decoder understands image embeddings and generates caption for text.

**Epoch Aware Masking** We first perform an simple token masking approach where we replace a random subset of words with a \<mask\> tokens each. It is similar to the BART Token Masking procedure during pre-training which we use to fine tune our model for image caption. We then increase the probability of masking as time progress in order to reduce the dependency of the decoder on the text and increase its dependency on images. We see that this model is dependent on the length of the caption being fed by the encoder and hence we provide two types of results, one where the model is fed masks with known caption length, and one where the model is fed masks with the length being the median length of all captions in the dataset.

**Text infilling** We also perform text infilling procedure outlined in the BART paper (Lewis et al., 2019). Random continuous words are replaced in the sentence with \<mask\> and the model learns to fill the \<mask\> tokens with the appropriate phrase. Our hope is that this model then learns not to depend on the text encoder for telling the decoder about the details in the image.
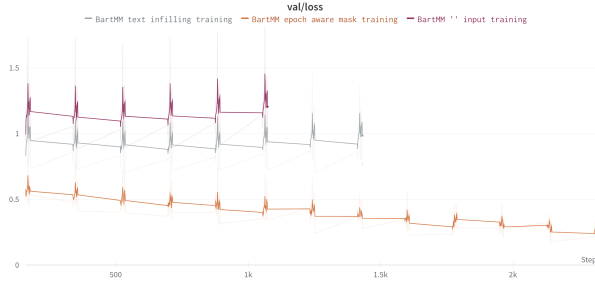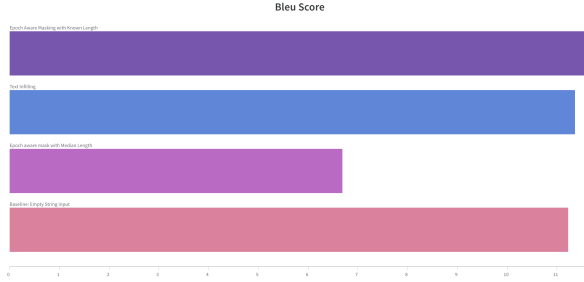
Figure 3: Caption Training Validation Loss Scores

| Method | BLEU | BERTScore |
|---|---|---|
| Baseline | 11.228 | - |
| Epoch Aware Mask (Length Known) | 11.662 | 0.6810 |
| Epoch Aware Mask (Median Length) | 6.683 | 0.5997 |
| Text Infilling | 11.365 | 0.7632 |

Table 2: Captioning performance



Figure 4: BLEU Scores



Figure 5: BERTScore

## 6 Results

We provide initial results based on validation loss, while training these three masking stratergies and use the following evaluation metrics.

- BLEU: Geometric mean of n-gram score with a brevity penalty to discourage shorter translation

- BERTScore: BERTScore is a similarity score for each token in the BERT representation of our output sequence with respect to each token the BERT representation of the reference ground truth sentence.

We also try to provide some initial reasoning behind them. From validation loss, we see the behaviour of our masking techniques. The epoch aware masking technique has a very small validation loss as the model is aware of the length of the caption it must generate. This dependency is all too common as we see when we compare BLEU and BERT scores. We see that due to text infilling, the model learns in a more robust fashion and hence has a lower validation loss when compared to the baseline empty string model.

**BLEU Score:** From figure 4 we see that the Epoch Aware model with known length input fairs the best out of all. This is expected but is one form of data leakage. The Text infilling model performs the next showing the importance of text infilling as
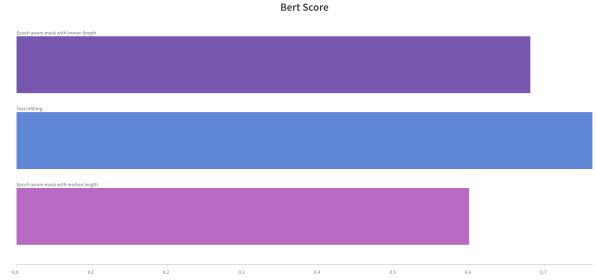
a masking scheme in BART pretraining. One more interesting observations we see is that the empty string baseline model is able to outperform the Epoch aware masking with median length, demonstrating that the Epoch aware masking model is heavily dependent on the length of the caption and has forgotten all about text infilling from its pretraining step. Something we hope to address during the course of our experiments.

**BERT Score:** In the interest of saving computational resources, we did not calculate BERT Score of our baseline. We see one important factor to note here. The BERT Score of the text infilling model is higher than both Epoch Aware masking techniques, demonstrating the better captioning capability of the text infilling model. This also suggests that the reason the BLEU score of Epoch aware masking model with known length is higher is mostly due to the length of the expected caption being known by the model due to which it can reduce the brevity penalty.

## 7 VQA Performance

We see the performance as shown in figures 6a, 6c, 6b. The first bar in each graph corresponds to the baseline VQA model with no pretraining. The second bar in each graph corresponds to the model being pretrained on the empty string captioning model. The third bar in each graph corresponds to the model being pretrained on the task infilling based captioning model. We see that the models
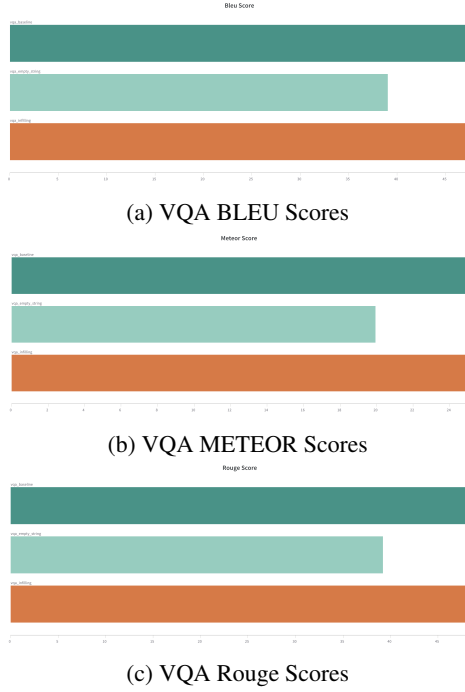
(a) VQA BLEU Scores



(b) VQA METEOR Scores



(c) VQA Rouge Scores

Figure 6: VQA Performances measured across BLEU, METEOR and ROUGE scores.

| Pretraining | Epoch | METEOR | Rouge | BLEU |
|---|---|---|---|---|
| Baseline | First Epoch | 24.00 | 46.73 | 45.75 |
| Baseline | Best Epoch | 24.90 | 48.24 | 47.42 |
| Empty String | Best Epoch | 19.90 | 39.21 | 39.03 |
| Infilling | First Epoch | 23.54 | 46.1 | 46.18 |
| Infilling | Best Epoch | 24.92 | 48.24 | 47.33 |

Table 3: VQA Performance

do not perform as expected. We believe that the models forget captioning right from the first epoch. We can see this from the table 3. In the interest of saving computational resources, we did not make an inference on the empty string model for the first epoch.

## 8  Error Analysis

We see that the BLEU Scores are not necessarily fair way to report these errors in generation.

We see some hallucination like in the first case 7a as London 2012 cannot be inferred from the example. In the second case, we opine that the model still is doing well, however its not very specific. As the expected captions all have wheel or tire in them, the model is unfairly penalised as having a bad BLEU score. This is an error in our evaluation scheme.



(a) Swimmers competing in the London 2012 Olympics. BLEU: 1.053e-229 BERT Score: 0.4122



(b) A man fixing a car. BLEU: 3.156e-76 BERT Score: 0.858

Figure 7: Examples of the type of captions made by text infilling
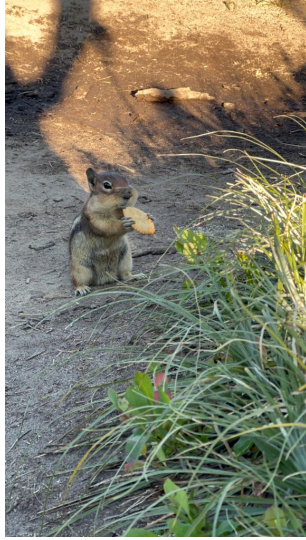
### 8.1  Prompting

We qualitatively evaluate captioning while being prompted. For the qualitative evaluations, we feed in prompts which are rather too informative, something we did to understand how the model behaves.

For image 1: 8a we prompt it as follows:

- Empty string model with no input
  Output: A man in black shirt and black pants looking outside the window of a dark coloured car

- Empty string model prompt: "<mask> <mask> <mask> <mask> gun <mask>"
  Output: A man in black shirt and black pants looking outside the window of a dark coloured car

- Epoch Aware model with prompt "<mask> <mask> <mask> <mask> <mask> <mask>"
  Output: A man in black shirt and black pants looking outside the window of a dark coloured car

- Empty string model with prompt: "<mask>

(a) Real World Example 1



(b) Real World Example 2

Figure 8: Evlauating Prompting characterstics

&lt;mask&gt; &lt;mask&gt; &lt;mask&gt; gun &lt;mask&gt;"
Output: A man shooting a gun .

- Text Infilling model with prompt: "&lt;mask&gt;"
  Output: A man looking through a window .

- Text Infilling model with prompt: "&lt;mask&gt;
  gun &lt;mask&gt;"
  Output: A man is getting ready to shoot a gun
  .

For image 2: 8b we prompt it as follows:

- Empty string model with no input
  Output: 'A black dog is climbing on a low
  tree branch.'

- Empty string model prompt: " A squirrel
  &lt;mask&gt;"
  Output: 'A black dog is climbing on a low
  tree branch.'

- Epoch Aware model with prompt "&lt;mask&gt;
  &lt;mask&gt;  &lt;mask&gt;  &lt;mask&gt;  &lt;mask&gt;
  &lt;mask&gt;"
  Output: A dog climbing a tree .

- Empty string model with prompt: "A squirrel
  &lt;mask&gt; &lt;mask&gt; &lt;mask&gt;"
  Output: A squirrel biting a tree

- Text Infilling model with prompt: "&lt;mask&gt;"
  Output: A black dog climbs a tree .

- Text Infilling model with prompt: "&lt;mask&gt;
  squirrel &lt;mask&gt;"
  Output: A black dog climbing a tree with a
  red squirrel

## 8.2 VQA

We notice that the 3 models fine-tuned on VQA
perform similarly but do not provide similar
answers.



Figure 9: The image is of a bicycle stand with bicycles
parked. The input question to the model is 'What is on
the ground?

We analyze the image shown in fig. 9. The im-
age is of a bicycle stand with bicycles parked. It
also has a red coloured building the background.
The input text to the model is 'What is on the
ground?'. The reference answers are:

- bikes

- bicycles

- bicycle

- bricks

- bike and bike racks

The outputs of our models are as follows:

- Baseline Model fine-tuned on VQA without pretraining for Image Captioning: bicycle

- Caption Empty string pretrained model fine-tuned on VQA: bike

- Caption Infilling pretrained model fine-tuned on VQA: bicycle

The model is successfully able to understand the intent of the question. It is also able to recognize the objects corresponding to the question in the image and generate an answer.

We notice that the models are correctly able to recognize the format of the input questions in most scenarios. However, it fails to do so in some scenarios.



Figure 10: An animated image in the COCO dataset. The input question is 'Why is the little girl's fingers dark?'

For image shown in fig. 10. The input text is 'Why is the little girl's fingers dark?' The reference answers are:

- there dirty

- shadow

- it's virtual realty

- because they are casting shadow in rendered game

The outputs of our models are as follows:

- Baseline Model fine-tuned on VQA without pretraining for Image Captioning: yes

- Caption Empty string pretrained model fine-tuned on VQA: woman

- Caption Infilling pretrained model fine-tuned on VQA: pink

The model fails to understand the intent of the question and hallucinates answers based on the image. The caption pretrained models recognize woman and the pink colour in the background and behaves as an object recognition model whereas the baseline model generates a 'yes'.



Figure 11: An image of landscape of a city in the COCO dataset. The input question is 'Where does the arrow point?'

For image shown in fig. 11, with given input text as 'Where does the arrow point?', the models are correctly able to understand the intent of the question but fail to answer correctly. The outputs of our models are as follows:

- Baseline Model fine-tuned on VQA without pretraining for Image Captioning: left

- Caption Empty string pretrained model fine-tuned on VQA: left

- Caption Infilling pretrained model fine-tuned on VQA: down

The models are correctly able to recognize the intent of questions that have either a 'yes' or 'no'

as answer. However, we do not analyze questions with 'yes-no' answers.



Figure 12: An image of a painted wall with graffiti in the COCO dataset. The input question is 'What color is the wall?'

For image shown in fig. 12, and input text 'What color is the wall?', the model is correctly able to understand the intent of the question The outputs of our models are as follows:

- Baseline Model fine-tuned on VQA without pretraining for Image Captioning: white

- Caption Empty string pretrained model fine-tuned on VQA: yellow

- Caption Infilling pretrained model fine-tuned on VQA: yellow

The wall in the input image contains multiple colours. All the models only generates one color as the output and fail to recognize the different colors in the image.

## 9   Conclusions and Future Work

From our experiments we have seen the utility of the the proposed approach in the captioning performance. We show that using text infilling based masking techniques results in a boost in performance over the baseline and other captioing techniques. However we are unable to show the utility of the proposed approach as a pretraining regime. We show from our error analysis that some of the errors of our models were due to unfair metrics. Hence we believe that there is a need for more research into better evaluation metrics. We also believe that the model is forgetful about its pretraining

regime when being finetuned on VQA. Hence the utility of pretraining could be focused on a Few Shot setting of VQA and other paired vision language understanding tasks.

## 10   Contributions

- Mukesh: He worked on Model Architecture, Token Infilling, Summarising model resuts and report, VQA Dataset prepration.

- Yogesh: He worked on Model Architecture, Baseline, Epoch Aware Masking, Captioing Dataset processing. VQA Training.

## References

Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2017. Bottom-up and top-down attention for image captioning and VQA. *CoRR*, abs/1707.07998.

Long Chen, Hanwang Zhang, Jun Xiao, Liqiang Nie, Jian Shao, and Tat-Seng Chua. 2016. SCA-CNN: spatial and channel-wise attention in convolutional networks for image captioning. *CoRR*, abs/1611.05594.

Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C. Lawrence Zitnick. 2015. Microsoft COCO captions: Data collection and evaluation server. *CoRR*, abs/1504.00325.

Xinlei Chen and C. Lawrence Zitnick. 2014. Learning a recurrent visual representation for image caption generation. *CoRR*, abs/1411.5654.

Xinpeng Chen, Lin Ma, Wenhao Jiang, Jian Yao, and Wei Liu. 2018. Regularizing rnns for caption generation by reconstructing the past with the present. *CoRR*, abs/1803.11439.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *CoRR*, abs/2010.11929.

Hao Fang, Saurabh Gupta, Forrest N. Iandola, Rupesh Kumar Srivastava, Li Deng, Piotr Dollár, Jianfeng Gao, Xiaodong He, Margaret Mitchell, John C. Platt, C. Lawrence Zitnick, and Geoffrey Zweig. 2014. From captions to visual concepts and back. *CoRR*, abs/1411.4952.

Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2019. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *CoRR*, abs/1910.13461.

Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. *CoRR*, abs/2004.06165.

Ron Mokady, Amir Hertz, and Amit H. Bermano. 2021. Clipcap: CLIP prefix for image captioning. *CoRR*, abs/2111.09734.

Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. 2015. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. *CoRR*, abs/1505.04870.

Matteo Stefanini, Marcella Cornia, Lorenzo Baraldi, Silvia Cascianelli, Giuseppe Fiameni, and Rita Cucchiara. 2021. From show to tell: A survey on image captioning. *CoRR*, abs/2107.06912.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *CoRR*, abs/1706.03762.

Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2014. Show and tell: A neural image caption generator. *CoRR*, abs/1411.4555.

Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. 2020. Unified vision-language pre-training for image captioning and vqa. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07):13041–13049.