

PART OF SPEECH MASKING EFFECT ON VISION-LANGUAGE REPRESENTATION LEARNING

by

Pasit Tiwawongrut

A Thesis Submitted in Partial Fulfillment of the Requirements for the Degree of
Master of Science in Data Science and Artificial Intelligence

Examination Committee: Dr. Chaklam Silpasuwanchai (Chairperson)
Dr. Chantri Polprasert
Dr. Attaphongse Taparugssanagorn

Nationality: Thai
Previous Degree: Bachelor of Computer Engineering
Khon Kaen University
Thailand

Scholarship Donor: Asian Institute of Technology

Asian Institute of Technology
School of Engineering and Technology
Thailand
July 2025

AUTHOR'S DECLARATION

I, Pasit Tiwawongrut, declare that the research work carried out for this thesis was in accordance with the regulations of the Asian Institute of Technology. The work presented in it are my own and has been generated by me as the result of my own original research, and if external sources were used, such sources have been cited. It is original and has not been submitted to any other institution to obtain another degree or qualification. This is a true copy of the thesis, including final revisions.

Date: 22 July 2025

Name: PASIT TIWAWONGRUT

Signature:

ACKNOWLEDGMENTS

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Aliquam malesuada bibendum arcu vitae elementum. Aliquet enim tortor at auctor. Sodales ut eu sem integer vitae justo. Tellus pellentesque eu tincidunt tortor aliquam nulla facilisi. Ut lectus arcu bibendum at varius vel pharetra. Integer vitae justo eget magna. Mollis nunc sed id semper risus. Donec ac odio tempor orci dapibus ultrices in iaculis. Non arcu risus quis varius quam quisque id diam vel. Felis donec et odio pellentesque diam volutpat commodo sed egestas. Dignissim diam quis enim lobortis scelerisque. Placerat vestibulum lectus mauris ultrices eros in cursus.

ABSTRACT

Vision language (VL) models have shown promising performance across multiple tasks in both zero-shot and fine-tuning setups. Most studies use masked language modeling as a pre-training task, applying random masking to image caption tokens. However, random token masking is not an optimal strategy for training VL models, and effective masking strategies in VL remain underexplored. In this work, we investigate the effects of part of speech (POS) masking, as each POS category contributes differently to sentence meaning. By pre-training models with different POS masking strategies, we evaluate each model on image-text retrieval and visual question answering tasks, categorizing each question type following the VALSE. Our findings contribute to a deeper understanding of how POS masking influences model performance, providing insights that can lead to more effective pre-training strategies for future VL models.

CONTENTS

	Page
TITLE PAGE	i
AUTHOR'S DECLARATION	ii
ACKNOWLEDGMENTS	iii
ABSTRACT	iv
LIST OF TABLES	vii
LIST OF FIGURES	viii
CHAPTER 1 INTRODUCTION	1
1.1 Background	1
1.2 Objective	2
1.3 Scope	3
CHAPTER 2 LITERATURE REVIEW	4
2.1 Vision-Language Representation Learning	4
2.2 Masked Language Modelling	5
CHAPTER 3 METHODOLOGY	7
3.1 Model architecture	7
3.2 Pre-training Objectives	8
3.2.1 Mask Language Modelling	8
3.2.2 Image-Text Contrastive Learning	9
3.2.3 Image-Text Matching	9
3.3 Part Of Speech Masking	10
3.4 Pre-Training Dataset	11
3.5 Evaluation	11
3.5.1 Image-Text Retrieval	11
3.5.2 Image-Text Matching	11
3.5.3 Visual Question Answering	12
3.6 Training	13

CHAPTER 4	Results	14
4.1	Pre-training	14
4.1.1	Flickr30K	14
4.1.2	POS Distribution in the Dataset	16
4.2	Image-Text Matching	17
4.2.1	VALSE	17
4.3	Visual question answering	18
4.4	Masking Ratio	18
CHAPTER 5	DISCUSSION	21
CHAPTER 6	CONCLUSION	22
6.1	Conclusion	22
6.2	Future work	22
REFERENCES		23

LIST OF TABLES

Tables	Page
Table 4.1 Flickr30K benchmark image retrieval result.	15
Table 4.2 VALSE benchmark for image-text matching result.	19
Table 4.3 VQA2.0 test-dev benchmark result.	19

LIST OF FIGURES

Figures	Page
Figure 3.1 Overall methodology	7
Figure 3.2 Visual question answering model architecture	12
Figure 4.1 MLM loss curves for different POS masking strategies (log scale).	16
Figure 4.2 ITC loss curves for different POS masking strategies (log scale).	17
Figure 4.3 ITM loss curves for different POS masking strategies (log scale).	18
Figure 4.4 Histogram of POS tag frequencies in the training dataset (sorted by frequency).	20

CHAPTER 1

INTRODUCTION

1.1 Background

Vision language (VL) models have gained significant attention due to their ability to perform both zero-shot and transfer learning, achieving high performance across numerous downstream tasks through pre-training with web-scale image-text pairs (Mo, Kim, Lee, & Shin, 2024; Z. Wang, Wu, Agarwal, & Sun, 2022; J. Zhang, Huang, Jin, & Lu, 2024). Many VL models incorporated masked language modeling (MLM) as a pre-training task, making it an important method to train VL models (J. Li et al., 2021; C. Li et al., 2022; Chen et al., 2020; W. Wang et al., 2023; Tan & Bansal, 2019). Typically, a subset of word tokens is randomly masked at a percentage during training, and the model is tasked with predicting these masked tokens using information from both visual and language modalities. This masking approach has proven to enhance the alignment between visual and linguistic representations, boosting performance in VL tasks (Tan & Bansal, 2019).

Despite the widespread adoption of MLM in VL training, its effects on model performance, efficiency, and training loss remain underexplored. Bitton, Stanovsky, Elhadad, and Schwartz (2021) demonstrated that many of the randomly masked tokens are often stop-words or punctuation, which the model can easily learn without any need for masking. Another study by Wilf et al. (2023) demonstrated that selectively masking infrequent words from the pre-training dataset can boost model performance on out-of-domain datasets during continued training. Additionally, Tou and Sun (2024) suggested that random masking causes the model to rely heavily on local text signals, and it results in inefficient and inconsistent interactions between modalities, leading to suboptimal performance. These find-

ings emphasize the importance of strategic token selection in MLM to enhance VL model performance and efficiency.

In this work, we aimed to address the gap in understanding how masking each part of speech (POS) impacts VL models. Each POS contributes distinctively to sentence meaning: nouns typically denote objects, while verbs describe actions and often demand contextual comprehension. By selectively masking different parts of speech, we could better understand how each POS category affects the alignment between visual and linguistic information. To further explore the effect of each POS, training without the MLM task and with different POS masking ratios were compared. The experiment is designed to answer the following questions:

1. How does masking each POS impact the performance, and training loss of VL pre-training models?
2. How does each POS masking strategy affect visual question answering (VQA) performance when analysed based on different question types?
3. What is the difference between training without the MLM task compared to training with it, and when masking each POS with a 100 percent masking ratio?

1.2 Objective

The objectives for our experiment are as listed.

1. Develop a pre-trained VL model to evaluate the impact of masking each POS on performance and training efficiency.
2. Benchmark the performance of each POS masking approach using specialized datasets to gain a deeper understanding of masking effects with retrieval and question answering tasks.
3. Train the model without POS masking and with POS masking at a 100 percent masking ratio.

1.3 Scope

1. The training and testing datasets are web-scale image-text pairs.
2. The model architecture is a cross-attention model, chosen for its ability to jointly predict answers based on information from multiple modalities.

CHAPTER 2

LITERATURE REVIEW

This section of the literature review is organized around two key topics relevant to our study. The first topic addresses VL models, providing an overview of the model architectures recently used in VL models and discussing the choice of the base architecture for the VL model used in this research. The second topic is MLM, an important pre-training approach that has improved VL model performance. Together, these sections provide a comprehensive overview of the methodological foundations of this study.

2.1 Vision-Language Representation Learning

VL learning aims to align visual and linguistic information for multimodal tasks that require reasoning across both modalities, such as image captioning, visual question answering, and multimodal retrieval. The training objective can be divided roughly into three main categories: contrastive, generative, and alignment. Firstly, the contrastive learning objective trains VL representations by maximizing the alignment score between paired images and text while minimizing the score between unpaired images and text (Radford et al., 2021; Jia et al., 2021; J. Yang et al., 2022). Secondly, the generative learning objective focuses on reconstructing masked tokens in either the image or text modality, or both, to learn VL representations (Singh et al., 2022; J. Li et al., 2021; Alayrac et al., 2022). This objective requires the model to utilize both modalities to reconstruct missing tokens, which enhances alignment. Lastly, the alignment objective involves learning VL representations by predicting whether an image and text pair match (Bao et al., 2022). The combined use of these three training objectives has proven effective and is commonly applied across various pre-trained VL models.

Recent advancements in VL fusion methods can be roughly categorized into three main approaches. The first approach is a separate unimodal encoder for each modality, as seen in models like CLIP (Radford et al., 2021) and ALIGN (Jia et al., 2021). This method is trained with the objective of aligning the intermediate outputs of each modality’s encoding. The second method uses a cross-attention layer to fuse multimodal inputs, e.g., Flamingo (Alayrac et al., 2022), LXMERT (Tan & Bansal, 2019), and ALBEF (J. Li et al., 2021). The cross-attention layer enables the model to fuse each modality more deeply. Finally, the third approach uses a single large attention model with concatenated image and text tokens as input, as in BEIT-3 (W. Wang et al., 2023), OSCAR (X. Li et al., 2020), UNITER (Chen et al., 2020), FLAVA (Singh et al., 2022), and mPLUG (C. Li et al., 2022). This approach allows for early-stage fusion of each modality, though it requires the highest amount of computational resources. In this work, we adopt the cross-attention method as the base model due to its effectiveness in fusing multimodal inputs. Additionally, this approach allows the model to be trained using the MLM task. We also use all three training objectives with a modified MLM for this experiment.

2.2 Masked Language Modelling

MLM is a widely used pre-training method in language model (LM) training (Devlin, Chang, Lee, & Toutanova, 2018; Lan, 2019; Yu et al., 2022; S. Zhang et al., 2022; Guu, Lee, Tung, Pasupat, & Chang, 2020) as a self-supervised task. BERT (Devlin et al., 2018) proposed MLM as a pre-training task, which has been proven effective for pre-training language models. The MLM task involves replacing some input tokens with a special [MASK] token, and the model must predict the masked tokens based on the given unmasked tokens. In the field of VL models, many VL models have also adopted MLM as a training task to train the model to predict masked text based on visual information (J. Li et al., 2021; C. Li

et al., 2022; Chen et al., 2020; W. Wang et al., 2023).

In the field of selective masking strategies in natural language processing, several works have further refined MLM to enhance training efficiency. ERNIE (Sun et al., 2019), SpanBERT (Joshi et al., 2020), and n -gram Masking (Levine et al., 2021) propose span masking instead of single-token masking, which forces the model to rely more on long-range dependencies rather than adjacent tokens, resulting in better performance compared to BERT (Devlin et al., 2018). Considering linguistic features, D. Yang, Zhang, and Zhao (2023) conducted a training analysis based on POS masking focused on LM training. The results showed that focusing the masking of non-function words (ADJ, ADV, NOUN, PROPN, and VERB) in the later stages of training can encourage the LM model to develop a better contextual understanding.

For selective masking in VL training, Bitton et al. (2021) introduced an object token masking strategy, selectively masking object tokens in image captions and pre-training the model. This approach achieved superior performance compared to random masking. Another study by Wilf et al. (2023) showed that selectively masking infrequent words from the pre-training dataset during continued training enhances model performance on out-of-domain datasets. Additionally, (Tou & Sun, 2024) proposed a curriculum-based masking strategy in which a reinforcement learning agent dynamically selects masking spans based on cross-modal interactions. This method improved the model’s multimodalities understanding while reducing the dataset size needed for effective training. In this work, we conduct experiments to analyze the impact of each POS on results within a VL setting.

CHAPTER 3

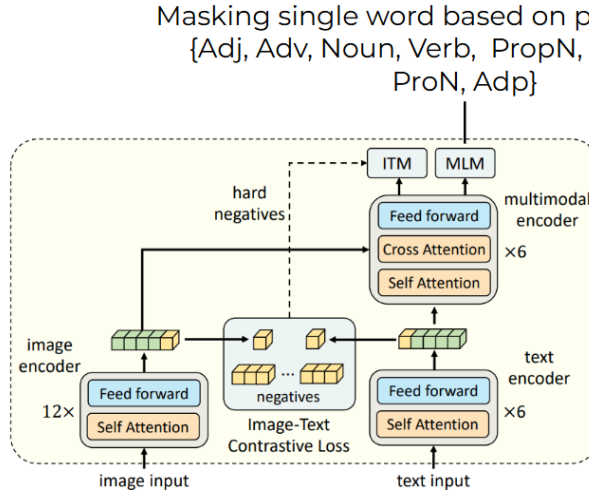
METHODOLOGY

In this chapter, the methodology is detailed as follows. First, we describe the architecture of the model. Second, we explain all pre-training loss functions used in this experiment. Third, the details of POS tagging are provided. Fourth, we outline the datasets used in this experiment. Lastly, we provide details on the visual question answering setup.

Figure 3.1

Overall methodology

Pre-training the model with a MLM task by masking tokens based on the POS in the image captions.



3.1 Model architecture

As shown in Figure 3.1, our model includes three main components: an image encoder, a text encoder, and a multimodal encoder. The first component is the

image encoder, for which we use ViT (Dosovitskiy et al., 2021), modified following (Radford et al., 2021), as the image encoder in this experiment. The second component is the text encoder, which employs a transformer architecture as BERT (Devlin et al., 2018) to encode image captions with BERT tokenizer for tokenization. The final component is the multimodal encoder, where VL interactions occur.

Given a training dataset D consisting of image-text pairs $(I_i, T_i) \in D$, where I_i is the image and T_i is the image caption of the i -th image, each image is first encoded as a sequence of tokens $\{v_{cls}, v_1, \dots, v_n\}$ using ViT (Dosovitskiy et al., 2021). Here, v_{cls} represents the embedding of the [CLS] token prepended to the image patch sequence. In this experiment, the image encoder was initialized with ViT-B-32 pre-trained on ImageNet-21K (Deng et al., 2009). Next, we use a 6-layer transformer, randomly initialized, to encode the image caption T_i into text embeddings $\{w_{cls}, w_1, \dots, w_n\}$, where w_{cls} is the embedding of the [CLS] token. Finally, both text and image encodings are passed through the multimodal encoder to fuse both inputs, producing multimodal encodings. For the multimodal encoder, a cross-attention layer is used, where both keys and values are the image encodings, and the text encoding serves as the query in the cross-attention layer.

3.2 Pre-training Objectives

In this work, we pre-train our model with three objectives: masked language modeling (MLM), image-text contrastive learning (ITC) and image-text matching (ITM).

3.2.1 Mask Language Modelling

Typically, a percentage of tokens $\{w_1, \dots, w_T\}$ are replaced with a special [MASK] token to create a masked caption T^{mask} . However, in this work, the masked tokens were selected based on POS type instead of randomly masking. The model

trained to predict the original tokens at the masked positions, conditioned on both the unmasked tokens in T^{mask} and the visual features of I as $p^{\text{mask}}(I, T^{\text{mask}})$. Let y^{mask} be a one-hot vector representing the ground-truth vocabulary for the masked token, where the masked token has a probability of 1. The model’s objective is to minimize the cross-entropy \mathbf{H} , given by:

$$\mathcal{L}_{\text{MLM}} = \mathbf{H}(y^{\text{mask}}, p^{\text{mask}}(I, T^{\text{mask}}))$$

For the masking ratio, each POS token is masked with either a 70 percent or 100 percent probability. In this work, random token masking was also tested with a masking ratio of 15 percent.

3.2.2 Image-Text Contrastive Learning

To improve each unimodal encoder’s representation, we used image-text contrastive learning to improve alignment of each modality. ITC aims to improve alignment by maximizing the similarity score of image and text from the same pair with the score function $s(I, T) = v_{cls}^\top w_{cls}$, and minimizing the similarity score of image and text not from its pair. We then calculate the softmax-normalized similarity score for each image to any text and each text to any image, identified as image-to-text $p^{i2t} \in \mathbb{R}^M$ and text-to-image $p^{t2i} \in \mathbb{R}^M$ scores as:

$$p_i^{i2t}(I) = \frac{\exp(s(I, T_i))/\tau}{\sum_{m=1}^M \exp(s(I, T_m)/\tau)}, \quad p_i^{t2i}(T) = \frac{\exp(s(T, I_i))/\tau}{\sum_{m=1}^M \exp(s(T, I_m)/\tau)}$$

where τ is a learnable temperature parameter. Let $y^{i2t}(I) \in \{0, 1\}^M$ and $y^{t2i}(T) \in \{0, 1\}^M$ be a ground truth with probability of 1 at a position of the same pair, and probability of 0 on the other hand. The ITC loss is calculated as cross-entropy \mathbf{H} between p and y :

$$\mathcal{L}_{\text{ITC}} = \frac{1}{2}(\mathbf{H}(y^{i2t}, p^{i2t}) + \mathbf{H}(y^{t2i}, p^{t2i}))$$

3.2.3 Image-Text Matching

To further improve multimodal alignment in the VL model, image-text matching was employed to enhance alignment. The model is trained to predict whether an

image and caption are from the same pair. A fully connected layer, followed by a softmax function, is added over the model. This layer takes the [CLS] embedding from the multimodal encoding as input to predict whether the pair is positive (matched) or negative (unmatched).

The loss function for ITM, using cross-entropy loss, is defined as:

$$\mathcal{L}_{\text{ITM}} = \mathbf{H}(y^{\text{itm}}, p^{\text{itm}}(I, T)),$$

where y^{itm} is a one-hot ground-truth label, and $p^{\text{itm}}(I, T)$ is the predicted class probability.

The full pre-training objective of our work can be written as:

$$\mathcal{L} = \mathcal{L}_{\text{MLM}} + \mathcal{L}_{\text{ITC}} + \mathcal{L}_{\text{ITM}}$$

3.3 Part Of Speech Masking

For each image caption, each token was classified into POS categories for masking. We used POS-tagging tools SpaCy¹ to classify each word into POS classes based on the Universal POS tag set². In this work, we modified the BERT tokenizer to integrate with SpaCy by using the Tokenizations³ tool to align BERT token IDs with SpaCy tokens IDs.

In this experiment, we explored the effect of each POS on VL learning in terms of performance, and training loss. For the main experiment, we focused on the following POS categories: nouns, verbs, adjectives, adverbs, proper nouns, determiners, auxiliaries, pronouns, and adpositions, each with a 70 percent masking probability. Each POS is also categorized as either functional or non-functional, where functional words include determiners, auxiliaries, pronouns, and adpositions, and non-functional words include nouns, verbs, adjectives, adverbs, and

¹POS-tagging tool SpaCy: <https://spacy.io/>

²Universal POS tag set: <https://universaldependencies.org/u/pos/>

³Tokenizations alignment library tool: <https://github.com/explosion/tokenizations>

proper nouns, for evaluation purposes. For the 100 percent masking setting, where all tokens corresponding to a specific POS are masked, we conducted experiments on non-functional parts of speech, including nouns, verbs, adjectives, adverbs, and proper nouns.

3.4 Pre-Training Dataset

We pre-trained the model on the Conceptual Captions dataset (Sharma, Ding, Goodman, & Soricut, 2018) and the MSCOCO dataset, totaling 2.4 million image-text pairs. In Conceptual Captions dataset, an automated process was used to select, filter, and refine these image-caption pairs to ensure they are clear, informative, and suitable for effective model training.

3.5 Evaluation

In this work, we evaluated each model trained with different types of POS masking through image-text retrieval, image-text matching, and visual question answering tasks. Details of the evaluation methods and datasets used in these tasks are provided in this section.

3.5.1 Image-Text Retrieval

For the image-text retrieval, the model was tested by performing zero-shot evaluations on the Flickr30K (Plummer et al., 2015) dataset for both image retrieval (IR) and text retrieval (TR). The Flickr30K dataset is used to assess the model’s overall performance in retrieval tasks. This setup allowed us to analyze how different POS masking strategies affect the model’s retrieval performance and the alignment between visual and textual representations.

3.5.2 Image-Text Matching

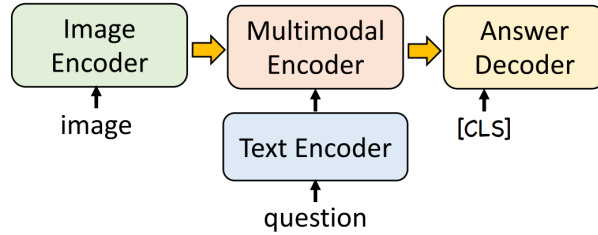
As demonstrated by Tou and Sun (2024), the results suggest that masking strategies impact a model’s ability to understand attributes, relationships, and word

order. In this work, we benchmarked each pre-trained model with specific POS masking against VALSE benchmark (Parcalabescu et al., 2022). For the VALSE dataset, this benchmark categorizes each image-text sample into different linguistic phenomena, including six distinct types: existence, plurality, counting, relation, action, and coreference. Each image caption in the VALSE dataset also includes a "Foil" version, where words related to each caption category are modified. This task is a classification task, where the model has to predict the correct caption for each image. We evaluated the model in a zero-shot manner by reusing the ITM head as a classifier. Evaluating models against this benchmark provides valuable insights into their semantic and contextual understanding of vision and language modality.

Figure 3.2

Visual question answering model architecture

Modified model architecture for VQA task.



3.5.3 Visual Question Answering

In this work, the visual question answering (VQA) task was treated as a classification task. A classification head was appended to generate the answer, as shown in Figure 3.2. The benchmark dataset for the VQA task is the VQA2.0 dataset (Goyal, Khot, Summers-Stay, Batra, & Parikh, 2017), which is constructed using images from COCO (Lin et al., 2014). This dataset includes 83,000 images for training, 41,000 for validation, and 81,000 for testing. We further train our model using the VQA2.0 training set.

3.6 Training

The model was pre-trained on a machine equipped with four NVIDIA A100 GPUs. The pre-training of the model was conducted using a batch size of 64 with 10 epochs. We used the AdamW optimizer with an initial learning rate of 1×10^{-4} and a weight decay of 0.02 to help regularize the training process. A cosine learning rate scheduler was applied, with the learning rate gradually increasing from a warm-up value of 1×10^{-5} during the first 5 epochs, before decaying towards a minimum learning rate of 1×10^{-5} by the end of training.

For the VQA task, the model was trained with a batch size of 32 on the same machine as pre-training. We used the AdamW optimizer with a learning rate of 2×10^{-5} and a weight decay of 0.02. A cosine learning rate scheduler was applied over 8 epochs, with a warm-up phase of 1 epoch starting at a learning rate of 1×10^{-5} , and decaying to a minimum of 1×10^{-6} .

CHAPTER 4

Results

This chapter presents all the experiment results for each experiment and evaluation, aimed at addressing the research questions. The results are divided into four sections, including pre-training, image-text matching, visual question answering, and the evaluation result of varying the POS masking percentage.

1. How does masking each POS impact the performance, and training loss of VL pre-training models?
2. How does each POS masking strategy affect visual question answering (VQA) performance when analysed based on different question types?
3. What is the difference between training without the MLM task compared to training with it, and when masking each POS with a 100 percent masking ratio?

4.1 Pre-training

To address the question of how masking each POS affects the performance and training loss of vision-language (VL) pre-training models, we present all relevant results in this section. All losses, including MLM, ITC, and ITM, along with the Flickr30K evaluation results, are provided. The loss values are plotted on a logarithmic scale to visualize improvements over time across different POS masking strategies. The results from training the ALBEF model using the same dataset are also included for consistent comparison.

4.1.1 Flickr30K

The Flickr30K evaluation results are shown in Table 4.1, which presents the top-1, top-5, and top-10 retrieval scores for both image-to-text (i2t) and text-to-image

(t2i) tasks across different training methodologies. Among the non-functional group, masking verbs and proper nouns leads to the most significant performance degradation compared to the random masking baseline, while masking nouns yields the best performance. The model with determiner masking achieves the highest overall performance, even surpassing the ALBEF model that uses knowledge distillation.

Masking Method	Flickr30K					
	TR			IR		
	r@1	r@5	r@10	r@1	r@5	r@10
ALBEF	70.40	89.50	94.00	54.66	82.02	88.70
Random Masking	67.00	88.00	93.75	52.61	80.14	87.76
Noun	67.15	88.60	94.65	52.73	80.45	87.79
Verb	54.85	82.85	90.05	43.82	73.84	82.82
Non-function Adj	62.30	87.30	92.40	47.39	75.47	84.06
Adv	46.85	76.25	85.75	36.40	66.38	76.78
PropN	44.85	74.40	84.10	34.91	64.09	75.01
Det	71.05	92.00	95.30	56.01	81.93	88.59
Function Aux	52.10	79.60	88.20	41.13	70.92	80.68
ProN	51.45	78.80	87.10	39.97	69.58	79.32
Adp	65.05	88.25	93.40	51.19	78.83	85.15

Table 4.1
Flickr30K benchmark image retrieval result.

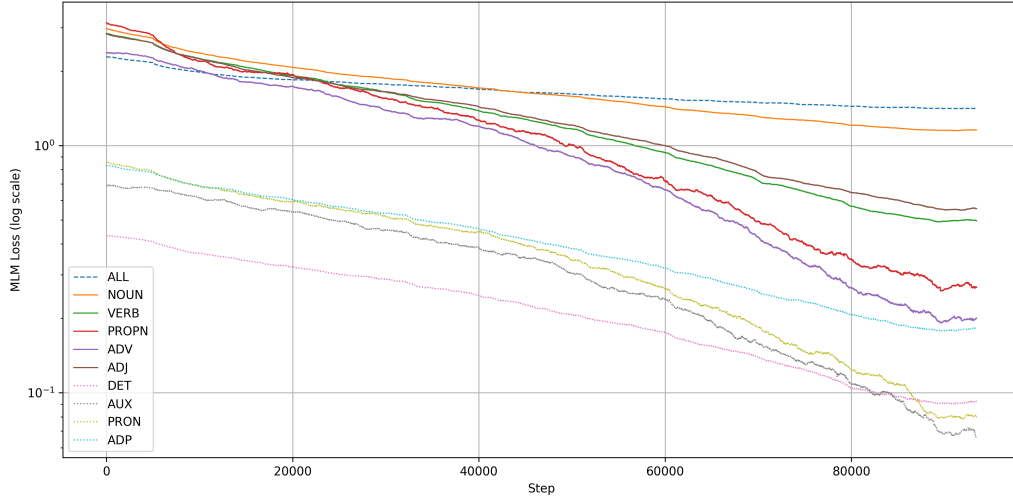
From the training loss curves, it is evident that different POS categories affect the convergence behavior in difference ways. The loss for MLM, ITC, and ITM are displayed in the Figure 4.1, Figure 4.2, and Figure 4.3 respectively For both ITM and ITC, the loss curves are similar in behavior and follow a consistent order relative to each other. In the MLM loss graph, we can see that POS masking in the functional group result in lower loss, while those in the non-functional group show higher loss, and the random masking show the highest loss by the end of training.

Taken together, the results show that masking each POS impacts both the training loss trajectory and final model performance in distinct ways. By observing the

MLM loss graph, we find that non-functional POS are more difficult for the model to learn through the MLM task, whereas functional POS are learned more quickly. By ranking the performance of each POS masking method, we observe that the ordering of POS aligns inversely with both the ITM and ITC loss curves. This indicates that the retrieval performance on the Flickr30K dataset is more strongly influenced by the ITM and ITC losses than by the MLM loss.

Figure 4.1

MLM loss curves for different POS masking strategies (log scale).

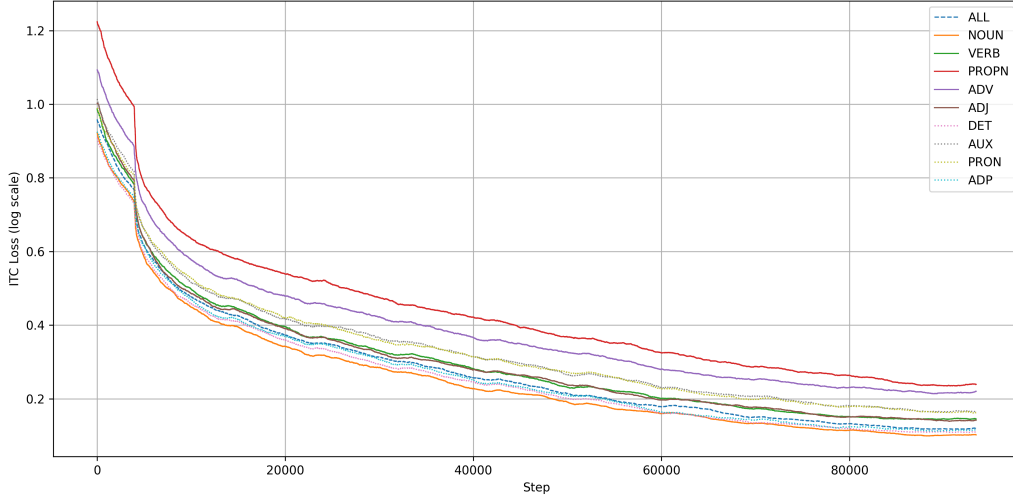


4.1.2 POS Distribution in the Dataset

To better understand the impact of POS masking, we analyzed the frequency distribution of each part-of-speech tag in the pre-training dataset. As shown in Figure 4.4, nouns are the most frequent category, followed by adpositions, determiners, and verbs. In contrast, categories such as interjections, symbols, and punctuation appear infrequently and were not considered in the masking strategies due to their limited semantic contribution.

Figure 4.2

ITC loss curves for different POS masking strategies (log scale).



4.2 Image-Text Matching

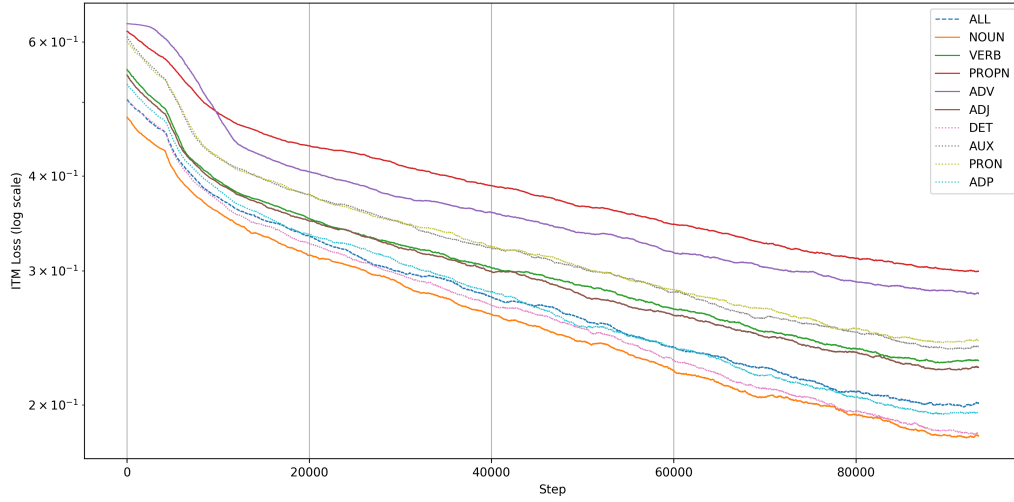
In this section, we evaluate the impact of POS masking on image-text matching by reporting zero-shot classification accuracy on the VALSE benchmark.

4.2.1 VALSE

Table 4.2 reports zero-shot classification accuracy on the VALSE benchmark for each POS masking strategy. The Random masking baseline yields a 60.20%, serving as a reference point. Masking nouns further raises the average to 60.45%, indicating that object words can be inferred from context in many cases. In contrast, masking verbs causes the average accuracy to drop to 58.66%, with the largest degradation observed on the Relation and Action categories. Adjective masking attains an average of 59.24%, showing that descriptive words contribute but are less critical than verbs. Masking adverbs and proper nouns leads to deeper declines, particularly on the Action and Coreference tasks, underscoring their role

Figure 4.3

ITM loss curves for different POS masking strategies (log scale).



in temporal and referential understanding.

These results highlight that the average accuracy (“Avg”) is the most indicative field for overall model competence, while the “Relation” and “Action” columns are especially informative for assessing the impact of verb and noun masking on fine-grained relational and dynamic reasoning.

4.3 Visual question answering

4.4 Masking Ratio

¹Spatial Relation

Masking Method	Existence quantifiers	Pularity number	Counting			VALSE		Action		Coreference		Foil-it!	Avg
			balanced	small number	adversarial	Sp.Re ¹ relations	replacement	actant swap	standard	clean			
Random Masking	65.06	61.43	54.64	57.81	62.83	61.61	68.04	51.88	49.70	43.37	85.79	60.20	
Non-function	Noun	67.63	62.60	52.59	54.64	64.39	59.84	68.15	48.87	51.31	49.21	85.69	60.45
	Verb	60.37	60.50	54.83	56.30	61.52	57.68	68.24	48.62	51.30	42.40	83.45	58.66
	Adj	60.85	60.55	54.00	56.84	67.01	57.68	65.68	50.92	50.34	44.74	83.01	59.24
	Adv	62.56	58.74	53.08	57.32	59.92	58.10	65.74	49.11	49.04	41.30	84.28	58.11
	PropN	61.51	59.23	52.49	56.25	61.26	55.86	64.31	50.85	50.36	43.03	82.62	57.98
Function	Det	60.14	63.33	53.47	57.86	65.40	59.06	66.67	50.43	50.09	38.99	87.94	59.40
	Aux	56.73	60.60	51.76	57.32	60.59	56.48	65.04	50.65	49.33	51.39	84.62	58.59
	ProN	56.05	61.33	50.39	54.88	58.87	58.93	64.36	48.05	53.23	50.48	83.40	58.18
	Adp	66.27	61.23	53.52	57.03	66.04	58.28	67.73	52.14	50.05	46.13	86.38	60.44

Table 4.2

VALSE benchmark for image-text matching result.

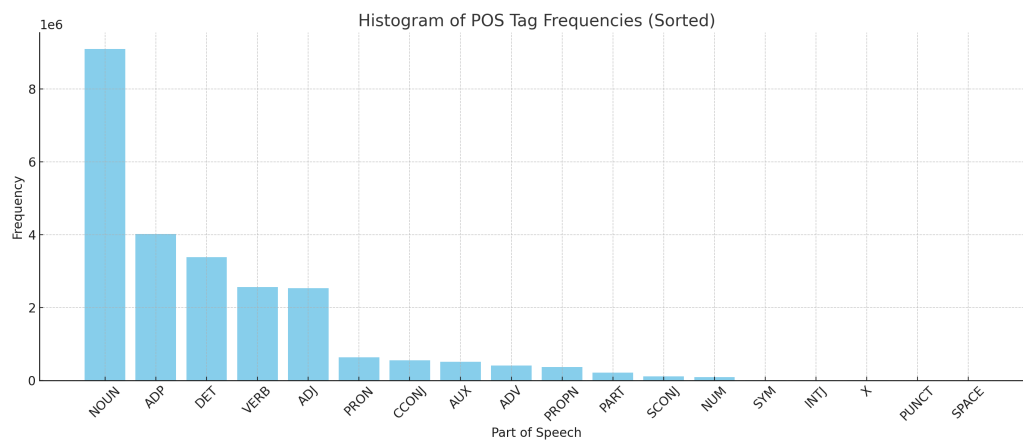
Masking Method		VQA2.0 test dev			
		Yes/No	Number	Other	Overall
Random Masking		87.88	49.64	59.63	70.28
Non-function	Noun	87.84	49.49	60.03	70.29
	Verb	87.17	48.39	58.43	69.13
	Adj	86.69	48.86	58.64	69.09
	Adv	83.10	43.83	52.49	64.12
	PropN	85.07	46.38	56.60	67.71
Function	Det	87.35	49.49	57.68	68.98
	Aux	85.25	46.59	56.24	67.09
	ProN	84.13	46.29	56.15	66.55
	Adp	87.07	48.82	58.07	68.96

Table 4.3

VQA2.0 test-dev benchmark result.

Figure 4.4

Histogram of POS tag frequencies in the training dataset (sorted by frequency).



CHAPTER 5

DISCUSSION

This suggests that verbs play a more critical role in aligning textual descriptions to visual features. Adjective masking results in a moderate performance drop, while other POS types (not shown here) follow similar trends. These findings reinforce the hypothesis that different parts of speech contribute unequally to cross-modal representation learning, with verbs carrying more alignment-critical information than nouns in this context.

CHAPTER 6

CONCLUSION

6.1 Conclusion

6.2 Future work

- Scalability for when combine with each pos - Compare with another method.

REFERENCES

- Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., ... others (2022). Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35, 23716–23736.
- Bao, H., Wang, W., Dong, L., Liu, Q., Mohammed, O. K., Aggarwal, K., ... Wei, F. (2022). Vlmo: Unified vision-language pre-training with mixture-of-modality-experts. *Advances in Neural Information Processing Systems*, 35, 32897–32912.
- Bitton, Y., Stanovsky, G., Elhadad, M., & Schwartz, R. (2021). Data efficient masked language modeling for vision and language. *arXiv preprint arXiv:2109.02040*.
- Chen, Y.-C., Li, L., Yu, L., El Kholy, A., Ahmed, F., Gan, Z., ... Liu, J. (2020). Uniter: Universal image-text representation learning. In *Computer vision – eccv 2020: 16th european conference, glasgow, uk, august 23–28, 2020, proceedings, part xxx* (p. 104–120). Berlin, Heidelberg: Springer-Verlag. Retrieved from https://doi.org/10.1007/978-3-030-58577-8_7 doi: 10.1007/978-3-030-58577-8_7
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition* (pp. 248–255).
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... Houlsby, N. (2021). An image is worth 16x16 words: Transformers for image recognition at scale. In *International conference on learning representations*. Retrieved from <https://openreview.net/forum?id=YicbFdNTTy>
- Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., & Parikh, D. (2017). Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 6904–6913).
- Guu, K., Lee, K., Tung, Z., Pasupat, P., & Chang, M.-W. (2020). *Realm: Retrieval-augmented language model pre-training*. Retrieved from <https://arxiv.org/abs/2002.08909>

- Jia, C., Yang, Y., Xia, Y., Chen, Y.-T., Parekh, Z., Pham, H., ... Duerig, T. (2021). Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning* (pp. 4904–4916).
- Joshi, M., Chen, D., Liu, Y., Weld, D. S., Zettlemoyer, L., & Levy, O. (2020). SpanBERT: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics*, 8, 64–77. Retrieved from <https://aclanthology.org/2020.tacl-1.5> doi: 10.1162/tacl_a.00300
- Lan, Z. (2019). Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*.
- Levine, Y., Lenz, B., Lieber, O., Abend, O., Leyton-Brown, K., Tennenholtz, M., & Shoham, Y. (2021). {PMI}-masking: Principled masking of correlated spans. In *International conference on learning representations*. Retrieved from <https://openreview.net/forum?id=3Aoft6NWFej>
- Li, C., Xu, H., Tian, J., Wang, W., Yan, M., Bi, B., ... others (2022). mplug: Effective and efficient vision-language learning by cross-modal skip-connections. *arXiv preprint arXiv:2205.12005*.
- Li, J., Selvaraju, R., Gotmare, A., Joty, S., Xiong, C., & Hoi, S. C. H. (2021). Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34, 9694–9705.
- Li, X., Yin, X., Li, C., Zhang, P., Hu, X., Zhang, L., ... others (2020). Oscar: Object-semantics aligned pre-training for vision-language tasks. In *Computer vision—eccv 2020: 16th european conference, glasgow, uk, august 23–28, 2020, proceedings, part xxx 16* (pp. 121–137).
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... Zitnick, C. L. (2014). Microsoft coco: Common objects in context. In *Computer vision—eccv 2014: 13th european conference, zurich, switzerland, september 6–12, 2014, proceedings, part v 13* (pp. 740–755).
- Mo, S., Kim, M., Lee, K., & Shin, J. (2024). S-clip: Semi-supervised vision-language learning using few specialist captions. *Advances in Neural Information Processing Systems*, 36.
- Parcalabescu, L., Cafagna, M., Muradjan, L., Frank, A., Calixto, I., & Gatt, A. (2022, May). VALSE: A task-independent benchmark for vision and language models centered on linguistic phenomena. In *Proceedings of the 60th annual meeting of the association for computational linguistics (volume 1: Long papers)* (pp. 8253–8280). Dublin, Ireland: Association for Compu-

- tational Linguistics. Retrieved from <https://aclanthology.org/2022.acl-long.567>
- Plummer, B. A., Wang, L., Cervantes, C. M., Caicedo, J. C., Hockenmaier, J., & Lazebnik, S. (2015). Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision* (pp. 2641–2649).
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... others (2021). Learning transferable visual models from natural language supervision. In *International conference on machine learning* (pp. 8748–8763).
- Sharma, P., Ding, N., Goodman, S., & Soricut, R. (2018). Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of acl*.
- Singh, A., Hu, R., Goswami, V., Couairon, G., Galuba, W., Rohrbach, M., & Kiela, D. (2022). Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 15638–15650).
- Sun, Y., Wang, S., Li, Y., Feng, S., Chen, X., Zhang, H., ... Wu, H. (2019). ERNIE: enhanced representation through knowledge integration. *CoRR*, *abs/1904.09223*.
- Tan, H., & Bansal, M. (2019). Lxmert: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 conference on empirical methods in natural language processing*.
- Tou, K., & Sun, Z. (2024, June). Curriculum masking in vision-language pre-training to maximize cross modal interaction. In K. Duh, H. Gomez, & S. Bethard (Eds.), *Proceedings of the 2024 conference of the north american chapter of the association for computational linguistics: Human language technologies (volume 1: Long papers)* (pp. 3672–3688). Mexico City, Mexico: Association for Computational Linguistics. Retrieved from <https://aclanthology.org/2024.naacl-long.203> doi: 10.18653/v1/2024.naacl-long.203
- Wang, W., Bao, H., Dong, L., Bjorck, J., Peng, Z., Liu, Q., ... others (2023). Image as a foreign language: Beit pretraining for vision and vision-language tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 19175–19186).
- Wang, Z., Wu, Z., Agarwal, D., & Sun, J. (2022). *Medclip: Contrastive learning from unpaired medical images and text*. Retrieved from <https://arxiv.org/abs/2210.10163>

- Wilf, A., Akter, S. N., Mathur, L., Liang, P. P., Mathew, S., Shou, M., ... Morency, L.-P. (2023). Difference-masking: Choosing what to mask in continued pretraining. *arXiv preprint arXiv:2305.14577*.
- Yang, D., Zhang, Z., & Zhao, H. (2023). *Learning better masking for better language model pre-training*.
- Yang, J., Li, C., Zhang, P., Xiao, B., Liu, C., Yuan, L., & Gao, J. (2022). Unified contrastive learning in image-text-label space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 19163–19173).
- Yu, W., Zhu, C., Fang, Y., Yu, D., Wang, S., Xu, Y., ... Jiang, M. (2022). *Dictbert: Enhancing language model pre-training with dictionary*. Retrieved from <https://arxiv.org/abs/2110.06490>
- Zhang, J., Huang, J., Jin, S., & Lu, S. (2024). Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(8), 5625-5644. doi: 10.1109/TPAMI.2024.3369699
- Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., ... Zettlemoyer, L. (2022). *Opt: Open pre-trained transformer language models*. Retrieved from <https://arxiv.org/abs/2205.01068>