



**NOVA**  
NOVA SCHOOL OF  
SCIENCE & TECHNOLOGY

DEPARTMENT OF  
COMPUTER SCIENCE

**AFONSO MIGUEL LOPES SIMPLÍCIO**

BSc in Mathematics

# **STRATEGIES TO BRIDGE MODALITIES IN LARGE VISION AND LANGUAGE MODELS**

MASTER IN ANALYSIS AND ENGINEERING OF BIG DATA

NOVA University Lisbon

September, 2024



**NOVA**

NOVA SCHOOL OF  
SCIENCE & TECHNOLOGY

DEPARTMENT OF  
COMPUTER SCIENCE

# STRATEGIES TO BRIDGE MODALITIES IN LARGE VISION AND LANGUAGE MODELS

**AFONSO MIGUEL LOPES SIMPLÍCIO**

BSc in Mathematics

**Adviser:** João Magalhães

*Full Professor, NOVA University Lisbon*

## **Examination Committee**

**Chair:** João Pires

*Associate Professor, NOVA University Lisbon*

**Rapporteur:** Bruno Martins

*Associate Professor, Instituto Superior Técnico*

**Adviser:** João Magalhães

*Full Professor, NOVA University Lisbon*

MASTER IN ANALYSIS AND ENGINEERING OF BIG DATA

NOVA University Lisbon

September, 2024

## **Strategies to Bridge Modalities in Large Vision and Language Models**

Copyright © Afonso Miguel Lopes Simplício, NOVA School of Science and Technology, NOVA University Lisbon.

The NOVA School of Science and Technology and the NOVA University Lisbon have the right, perpetual and without geographical boundaries, to file and publish this dissertation through printed copies reproduced on paper or on digital form, or by any other means known or that may be invented, and to disseminate through scientific repositories and admit its copying and distribution for non-commercial, educational or research purposes, as long as credit is given to the author and editor.

## ACKNOWLEDGEMENTS

First and foremost, I would like to express my deepest gratitude to my advisor, Professor João Magalhães, for granting me the opportunity to explore this exciting field of study. His unwavering guidance, encouragement, and insights were instrumental throughout the duration of my thesis, and for that, I am profoundly grateful.

I would also like to extend my thanks to NOVA School of Science and Technology for providing me with the tools, knowledge, and academic environment that have been critical in enabling me to complete this work.

To my friends and family, thank you for your constant encouragement and belief in me over the years. Your support has been a cornerstone of my journey. A special mention goes to my girlfriend for her steadfast patience and assistance throughout this challenging period.

# ABSTRACT

Vision and language models have achieved remarkable results by leveraging cross-modality encoders that connect visual and textual information. However, training such models with large amounts of data is computationally expensive. To address this, recent models have employed frozen pre-trained image encoders and large language models (LLMs), focusing on efficiently transforming visual embeddings into a format that can be processed by text decoders.

In this dissertation, we investigate the methods for establishing this vision-to-language connection. Building upon the framework used to train the Fromage [14] model, we conduct ablation studies to explore various techniques for extracting and transforming visual embeddings to be used as input of the text decoder. We focus on two primary architectures: (1) one that extracts a single visual embedding from the encoder and converts it into a sequence of embeddings using a linear transformation, and (2) one that extracts all token embeddings from a given layer and transforms them into a sequence of embeddings using cross-attention. Our experiments examine the impact of varying the number of embeddings generated and the choice of the encoder layer from which the embeddings are extracted. We evaluate model performance across several tasks, including image captioning, retrieval, and visual question answering. Based on our findings, we conclude that utilizing multiple visual token embeddings as decoder inputs significantly improves model performance.

Additionally, following this framework, we trained a new generative model for European Portuguese. As part of this effort, we created a Portuguese vision-and-language dataset by translating existing English datasets using machine translation.

**Keywords:** Visual Question Answering, Natural Language Processing, Multimodal Models, Contrastive Learning, Retrieval

## RESUMO

Modelos de visão e linguagem têm alcançado resultados notáveis ao utilizarem encoders multimodais que conectam informação visual e textual. No entanto, treinar esses modelos com grandes quantidades de dados é computacionalmente caro. Para contornar esta limitação, modelos recentes têm utilizado encoders de imagem pré-treinados e congelados, assim como Large Language Models (LLMs), concentrando-se na transformação eficiente de embeddings visuais para um formato que possa ser processado por decoders de texto.

Nesta dissertação, investigamos os métodos usados para estabelecer esta conexão entre visão e linguagem. Baseando-nos na arquitetura usada para treinar o modelo Fromage [14], realizamos estudos de ablação para explorar várias técnicas de extração e transformação de embeddings visuais para serem usados como entrada do decoder de texto. Focamos em duas arquiteturas principais: (1) uma que extrai um único embedding visual do encoder e converte numa sequência de embeddings através de uma transformação linear, e (2) uma que extrai todos os tokens de uma camada específica e transforma-os utilizando cross-attention.

As nossas experiências examinam o impacto de variar o número de embeddings gerados e a escolha da camada do encoder de onde são extraídos os embeddings. Avaliamos o desempenho do modelo em várias tarefas, incluindo legendagem de imagens, recuperação de imagens e resposta a perguntas visuais. Com base nos nossos resultados, concluímos que a utilização de vários embeddings visuais como entrada do decoder melhora significativamente o desempenho do modelo.

Além disso, seguindo este enquadramento, treinamos um novo modelo generativo para português europeu. Como parte deste esforço, criamos um conjunto de dados de visão e linguagem em português, traduzindo conjuntos de dados existentes em inglês utilizando tradução automática.

**Palavras-chave:** Visual Question Answering, Processamento de Linguagem Natural, Modelos Multimodais, Aprendizagem Contrastiva, Retrieval

# CONTENTS

<b>List of Figures</b>	<b>viii</b>
<b>List of Tables</b>	<b>xi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Context and Motivation . . . . .	1
1.2 Problem Definition and Objective . . . . .	1
1.3 Contributions . . . . .	3
1.4 Document Structure . . . . .	4
<b>2 Background</b>	<b>6</b>
2.1 Elements of Vision and Language Model Architectures . . . . .	6
2.1.1 Transformer . . . . .	6
2.1.2 Image Encoders . . . . .	9
2.2 Learning . . . . .	11
2.2.1 Loss functions . . . . .	11
2.2.2 Tasks . . . . .	15
2.2.3 Model Optimization . . . . .	16
2.2.4 LoRA . . . . .	20
2.2.5 Batch creation . . . . .	20
2.3 Summary . . . . .	21
<b>3 Related Work</b>	<b>23</b>
3.1 Vision and Language Model Encoder Architectures . . . . .	23
3.1.1 A Unified Mathematical Framework . . . . .	23
3.1.2 Single-stream attention . . . . .	24
3.1.3 Dual-stream attention . . . . .	25
3.1.4 Cross-attention . . . . .	26
3.1.5 Applications of V&L Encoder Models . . . . .	27
3.2 Vision and Language Model Decoder Architectures . . . . .	27

3.2.1	Large Language Models . . . . .	28
3.2.2	ViT Token extraction . . . . .	29
3.2.3	Bridging function . . . . .	29
3.2.4	Multimodal bridging function . . . . .	31
3.2.5	Text and Image in the Decoder . . . . .	32
3.2.6	Extra Vision-Language (VL) Features . . . . .	33
3.3	Summary . . . . .	35
<b>4</b>	<b>V&amp;L Tasks and Datasets</b>	<b>36</b>
4.1	Datasets . . . . .	36
4.1.1	MSCOCO . . . . .	36
4.1.2	VisDial . . . . .	37
4.1.3	CC3M . . . . .	37
4.2	Tasks . . . . .	38
4.2.1	Captioning . . . . .	38
4.2.2	Retrieval . . . . .	39
4.2.3	Visual Question Answering . . . . .	40
<b>5</b>	<b>A Portuguese Vision and Language Model</b>	<b>41</b>
5.1	Introduction . . . . .	41
5.2	EU-Portuguese Datasets for Vision and Language AI . . . . .	42
5.3	Method . . . . .	45
5.3.1	V-Glória Architecture . . . . .	45
5.3.2	Training . . . . .	46
5.3.3	Implementation details . . . . .	47
5.4	Experimental Setup . . . . .	47
5.5	Results and Discussion . . . . .	48
5.5.1	Cross-Modal Retrieval Results . . . . .	48
5.5.2	Image Captioning Results . . . . .	48
5.5.3	Visual Dialog Results . . . . .	49
5.6	Conclusions . . . . .	52
<b>6</b>	<b>Ablation Study: CLS based models</b>	<b>53</b>
6.1	ViT Layer . . . . .	54
6.2	Multiple Vision Tokens . . . . .	58
6.2.1	Layer $L$ . . . . .	59
6.2.2	Layer $L - 1$ . . . . .	62
6.3	Other CLS Based Approaches . . . . .	66
6.4	Summary . . . . .	71
<b>7</b>	<b>Ablation Study: All Tokens Based Models</b>	<b>72</b>
7.1	Methods . . . . .	72

7.2	Layer . . . . .	79
7.3	Number of tokens . . . . .	82
7.3.1	L . . . . .	83
7.3.2	L-1 . . . . .	87
7.4	Summary . . . . .	91
<b>8</b>	<b>Conclusions</b>	<b>92</b>
8.1	Contributions . . . . .	92
8.2	Take Away Messages . . . . .	92
8.3	Future work . . . . .	93
	<b>Bibliography</b>	<b>94</b>

## LIST OF FIGURES

2.1	Architecture of the Original Transformer, with an encoder and a decoder. [43].	6
2.2	Structure of a transformer decoder. [11]	9
2.3	Computation of a convolution [43].	10
2.4	Visual Transformer Architecture [10]	11
2.5	Illustration of the triplet loss. By minimizing the loss, the negative pair stays farther than the positive one from the anchor by the margin.	13
2.6	Comparison and examples of different scheduling functions. [43]	19
2.7	Example of a cosine scheduler function with a linear warmup. [43]	20
2.8	a) shows the nuclear norm of a noisy gradient estimate for various levels of noise (b) compares different sampling strategies. [41]	21
3.1	Visualization of the different score matrices [4]	24
3.2	Encoder architectures [4]	25
3.3	The single stream architecture of the ViLT model. [12]	25
3.4	The dual stream architecture of the CLIP model. [28]	26
3.5	The cross-attention architecture of the LXMERT model. Note that the initial layers are dual-stream. [32]	26
3.6	An evolutionary graph of the research work conducted on LLaMA. [46]	28
3.7	Architecture of the LLaVA model. The image here represented as $X_v$ is transformed to $Z_v$ by the vision Encoder entering the LLM as $H_v$ after the linear projection. [21]	30
3.8	Architecture and training of the Fromage Model. The output of the Visual encoders is transformed by the linear layer $W_c$ when entering the LLM and transformed by the linear layer $W_i$ when used for retrieval, being then compared with the transformation of the vector [RET] by the layer $W_i$ . [14]	31
3.9	InstructBLIP [8]	32

3.10	The architecture of CogVLM. (a) The illustration about the input, where an image is processed by a pre-trained ViT and mapped into the same space as the text features. (b) The Transformer block in the language model. The image features have a different QKV matrix and FFN. Only the purple parts are trainable. . . . .	33
4.1	Example images and captions from the Microsoft COCO Caption dataset. [18]	37
4.2	Example of a Visual Dialog [9] . . . . .	38
4.3	Example of images, alt-text and conceptual caption from the CC3M dataset [31].	38
5.1	Overview of the V-Glória architecture. The model is trained on image-text pairs for image captioning and image-text retrieval. The LLM and visual encoder are frozen, while the three projection layers (in yellow), with weight matrices $\mathbf{W}_c$ , $\mathbf{W}_i$ , and $\mathbf{W}_t$ , are learned. . . . .	46
5.2	V-Glória can solve core vision and language tasks. . . . .	51
6.1	Architecture of our model generative part where $h$ is a linear layer that transforms a CLS token, where we change the layer we extract the token. The solid black line corresponds to the layer $L$ used by the original model, but we try for the last 4 layers. . . . .	55
6.2	Architecture of our model generative part where $h$ transforms the CLS token of the last layer into $k$ different tokens, here we change the $k$ . . . . .	59
6.3	Architecture of our model generative part where $h$ transforms the CLS token from layer $L - 1$ into $k$ different tokens, here we change the $k$ . . . . .	62
6.4	Architecture of the generative part of the Last 4 Layers model, where we extract the CLS tokens of the last 4 layers and transform each one with a different linear transformation. . . . .	66
6.5	Architecture of the generative part of the Last 2 Layers model, where we extract the CLS tokens of the last 2 layers and transform each one with a different linear transformation. . . . .	67
6.6	Architecture of the generative part of the L2Lx2 model, where we extract the CLS tokens of the last 2 layers and transform each one with two different linear transformations creating 4 embeddings. . . . .	67
6.7	Architecture of the generative part of the CLS comparing model, where we use the CLS tokens from the last 2 layers and create 2 new embeddings by subtracting and multiplying element-wise these two. Obtaining 4 tokens where each one is transformed with a different liner transformation. . . . .	68
7.1	Architecture of the generative part of the Grouped Concatenation with Linear Transformation (Grouped) model, where we extract the tokens from layer $L$ except the CLS and group them in 4 groups where the vectors inside a group are concatenated, obtaining 4 vectors which are transformed by a linear layer.	73

7.2	Architecture of the generative part of the Position-Based Grouping with linear transformation (Regions) model, where we extract the tokens from layer $L$ except the CLS and group them in 5 regions corresponding to the region each token represents. The vectors inside a group are concatenated, obtaining 5 vectors which are transformed by a linear layer. . . . .	74
7.3	Architecture of the generative part of the weighted average (W avg) model, where we extract the tokens from layer $L$ except the CLS and group them in 4 groups where we obtain a weighted average of the embeddings in each group, obtaining 4 vectors which are transformed by a linear layer. . . . .	74
7.4	Architecture of the generative part of the Cross-Attention model, where we use 4 learned queries which interact with the image tokens from the last layer through cross-attention outputting 4 token embeddings. . . . .	75
7.5	Architecture of the generative part of the Cross-Attention model, where we choose from which layer to extract all tokens before they interact with the learned queries through cross-attention. The solid lines show layer $L$ which is the original layer we used the previous experiments, here we try different layers. . . . .	79
7.6	Architecture of the generative part of the Cross-Attention model, where we use $k$ learned queries which interact with the image tokens from the last layer through cross-attention outputting $k$ token embeddings, here we change the number $k$ . . . . .	83
7.7	Architecture of the generative part of the Cross-Attention model, where we use $k$ learned queries which interact with the image tokens from layer $L - 1$ through cross-attention outputting $k$ token embeddings, here we change the number $k$ . . . . .	87

## LIST OF TABLES

4.1	Comparison between the original number of images and the actually available images in the cc3m [31] dataset. . . . .	38
5.1	Translation results of sample captions from the CC3M dataset, using each of the three considered translation approaches. . . . .	44
5.2	Translation statistics, for CC3M and COCO, with different machine translation approaches. # Samples - total number of samples, # Tokens - total number of tokens, # Avg. Tokens/Sample - average number of tokens per sample. * Stands for the original captions. . . . .	44
5.3	Retrieval results for CC3M PT-PT and MSCOCO PT-PT datasets. . . . .	48
5.4	Captioning results on the validation split of the CC3M PT-PT and MSCOCO PT-PT datasets. . . . .	49
5.5	Zero-shot results on VisDial [9], for image-and-text-to-text (IT2T) and text-to-image (T2I) retrieval. Unlike previous methods, is capable of generating free-form text interleaved with image outputs through text-to-image retrieval. . . . .	49
6.1	CC3M validation metrics for different layers. . . . .	54
6.2	VQA validation results for different layers. . . . .	55
6.3	MSCOCO captioning validation metrics for different layers. . . . .	55
6.4	MSCOCO retrieval validation metrics for different layers. . . . .	55
6.5	VQA answers generated by the models extracting the cls embedding from different layers. . . . .	56
6.6	Examples of captions generated by the models extracting the cls embedding from different layers on images from the MSCOCO dataset. . . . .	57
6.7	CC3M validation metrics for different numbers of visual tokens with CLS extracted from the last layer . . . . .	59
6.8	VQA validation metrics for different numbers of visual tokens with CLS extracted from the last layer . . . . .	59
6.9	MSCOCO validation metrics for different numbers of tokens extracted from the last layer. . . . .	60

6.10	MSCOCO retrieval validation metrics for different numbers of token embeddings with CLS token extracted from the last layer. . . . .	60
6.11	VQA answers generated by the models with different numbers of tokens extracting the cls embedding from the last layer. . . . .	60
6.12	Examples of captions generated by the models with different number of visual tokens with CLS token extracted from the last layer on images from the MSCOCO dataset. . . . .	61
6.13	CC3M validation metrics for different numbers of visual tokens with CLS extracted from the second to last layer . . . . .	63
6.14	VQA validation metrics for different numbers of visual tokens with CLS extracted from the second to last layer. . . . .	63
6.15	MSCOCO validation metrics for different numbers of visual tokens with CLS extracted from the second to last layer. . . . .	63
6.16	MSCOCO retrieval validation metrics for different numbers of token embeddings with CLS token extracted from the second to last layer. . . . .	63
6.17	Examples of captions generated by the models with different number of visual tokens with CLS token extracted from the second to last layer on images from the MSCOCO dataset. . . . .	64
6.18	VQA answers generated by the models with different numbers of tokens extracting the cls embedding from the second to last layer. . . . .	65
6.19	CC3M validation metrics for different other CLS Based approaches . . . . .	68
6.20	VQA validation metrics for different other CLS Based approaches . . . . .	68
6.21	MSCOCO validation metrics for different other CLS Based approaches . . . . .	68
6.22	MSCOCO retrieval validation metrics for different other CLS Based approaches . . . . .	69
6.23	VQA answers generated by the models with different numbers of tokens extracting the cls embedding from the second to last layer. . . . .	69
6.24	Examples of captions generated by the different CLS based approaches on images from the MSCOCO dataset. . . . .	70
7.1	CC3M validation metrics for different methods using all tokens. . . . .	75
7.2	VQA validation metrics for different methods using all tokens. . . . .	76
7.3	MSCOCO captioning validation metrics for different methods using all tokens. . . . .	76
7.4	MSCOCO retrieval validation metrics for different methods using all tokens. . . . .	76
7.5	Examples of captions generated by different models using various methods to extract all tokens for images from the MSCOCO dataset. The empty captions of the Regions model correspond to images this model returned nothing while trying to caption. . . . .	77
7.6	VQA answers generated by the models using different methods to process all visual tokens. . . . .	78
7.7	CC3M validation metrics for different layers using cross-attention. . . . .	80
7.8	VQA validation metrics for different layers using cross-attention. . . . .	80

7.9	MSCOCO captioning validation metrics for different layers using cross-attention.	80
7.10	MSCOCO retrieval validation metrics for different layers using cross-attention.	80
7.11	Examples of captions generated by the models extracting the cls embedding from different layers on images from the MSCOCO dataset. . . . .	81
7.12	VQA answers generated by the models with different numbers of tokens extracting the cls embedding from the second to last layer. . . . .	82
7.13	CC3M validation metrics for different numbers of tokens using cross-attention extracting the tokens from the last layer. . . . .	83
7.14	VQA validation metrics for different numbers of tokens using cross-attention extracting the tokens from the last layer. . . . .	83
7.15	MSCOCO captioning validation metrics for different numbers of tokens using cross-attention extracting the tokens from the last layer. . . . .	84
7.16	MSCOCO retrieval validation metrics for different numbers of tokens using cross-attention extracting the tokens from the last layer. . . . .	84
7.17	Examples of captions generated by the Cross Attention models using different numbers of tokens and extracting all tokens from the last layer on images from the MSCOCO dataset. . . . .	85
7.18	VQA answers generated by the models with different numbers of tokens extracting the cls embedding from the second to last layer. . . . .	86
7.19	CC3M validation metrics for different numbers of tokens using cross-attention extracting the tokens from the second to last layer. . . . .	87
7.20	VQA validation metrics for different numbers of tokens using cross-attention extracting the tokens from the second to last layer. . . . .	87
7.21	MSCOCO captioning validation metrics for different numbers of tokens using cross-attention extracting the tokens from the second to last layer. . . . .	88
7.22	MSCOCO retrieval validation metrics for different numbers of tokens using cross-attention extracting the tokens from the second to last layer. . . . .	88
7.23	Examples of captions generated by the Cross Attention models using different numbers of tokens and extracting all tokens from the second to last layer on images from the MSCOCO dataset. . . . .	89
7.24	VQA answers generated by the Cross Attention models with different numbers of tokens extracting all tokens from the second to last layer. . . . .	90

# INTRODUCTION

## 1.1 Context and Motivation

Vision and language are two of the most extensively studied modalities in AI, with significant progress achieved across various tasks. Recently, there has been a tremendous advance in multimodal models, particularly those that combine vision and language. These models enable a more comprehensive understanding of the world by integrating visual and textual information, leading to breakthroughs in tasks such as image captioning, visual question answering, and image-text retrieval. Advancing these vision-and-language models is crucial for the continued evolution of AI and its ability to understand and interact with the world.

However, training large vision and language models from scratch can be prohibitively expensive. To address this, it is essential to explore methods for adapting existing models. This is often achieved by leveraging pre-trained vision models and pre-trained language models. A common approach is to use a vision model as an encoder, which transforms an image into a vector or a set of vectors, known as embeddings, that represent the image's content. Meanwhile, a large language model can be used as a decoder, capable of generating text from a text input and the image embeddings.

Combining vision and language models is not trivial, and there are numerous ways to design such systems. In this work, we focus on models where both the image encoder and the text decoder are frozen—their weights remain unchanged during training—and we only train a function that transforms the image embeddings. However, there are other approaches, such as training models where the encoder and decoder are not frozen, or using more complex methods to connect these two components.

## 1.2 Problem Definition and Objective

The general research problem of this dissertation lies in how to bridge the visual and the textual modalities, specifically how to bridge a visual encoder with a text decoder. To explain our problem, let's start by defining what we mean by a visual encoder and a text

decoder. We define these models as the following functions:

- **Vision Encoder:**

$$f(I) = Y_v, \quad (1.1)$$

where  $I$  is an image and  $Y_v \in \mathbb{R}^{N \times d_v}$  represents a sequence of  $N$  visual embeddings, each one being a vector of dimension  $d_v$ .

- **Text Decoder:**

$$g(Z) = T_{output}, \quad (1.2)$$

where  $Z \in \mathbb{R}^{M \times d_t}$  is a sequence of  $M$  text embeddings each one being a vector of dimension  $d_t$ , and  $T_{output} \in \mathbb{R}^{p \times d_t}$  is the generated sequence of  $p$  output text tokens that are then transformed into text.

This means that our visual encoder transforms an image into a sequence of vectors (or 1 vector if  $N = 1$ ), and the text decoder transforms a sequence of vectors (or 1 vector if  $M = 1$ ) into text. In a Generative Vision and Language model, the goal is for the decoder to generate text based on the image  $I$  as well as an input text  $T_{input}$ .

To enable the text decoder to generate text based on both the image and the input text, we need a way to map the visual information into the same space as the text tokens. Let  $Z_t$  be the sequence of embeddings corresponding to the input text  $T_{input}$ . We define a bridging function  $h$  that transforms the visual embeddings  $Y_v$  into a compatible format with the text embeddings. The output of this bridging function is  $Z_v$ , a sequence of vectors in the same dimensional space as the text embeddings:

$$h(Y_v) = Z_v \in \mathbb{R}^{M_v \times d_t} \quad (1.3)$$

Here,  $Z_v$  represents the visual embeddings in the text embedding space after the transformation, being composed of  $M_v$  embedding with dimension  $d_t$ . Note that our function  $h$  transforms  $N$  vectors into  $M_v$  vectors, these values may not be the same depending on what function we use, as it will become clearer when we present the different functions later on.

Once we have  $Z_v$ , the transformed visual embeddings, we combine them with the text embeddings  $Z_t$  from the input text  $T_{input}$ . The combined sequence  $Z$  is formed by concatenating the visual and text embeddings,

$$Z = [Z_v, Z_t] \in \mathbb{R}^{(M_v + M_t) \times d_t},$$

where the combined sequence  $Z$  is then passed to the text decoder  $g(Z)$ , which processes the mixed modality input and generates the output text  $T_{output}$ . In a practical example, consider an image  $I$  representing a dog in a park and an input prompt  $T_{input}$ , such as "A dog is". The process would proceed as follows:

1. **Visual Encoding:** The image  $I$  is passed through the visual encoder  $f(I)$ , producing visual embeddings  $Y_v$ .

2. **Text Encoding:** The input text  $T_{input}$  is converted into a sequence of text embeddings  $Z_t$ .
3. **Bridging Function:** The visual embeddings  $Y_v$  are transformed by the bridging function  $h(Y_v)$  into  $Z_v$ , so that they reside in the same embedding space as the text tokens.
4. **Concatenation:** The visual embeddings  $Z_v$  and text embeddings  $Z_t$  are concatenated into a unified sequence  $Z = [Z_v, Z_t]$ .
5. **Output Generation:** The combined sequence  $Z$  is passed to the text decoder, which generates the final output text, such as "A dog is running in the park".

The central objective of this dissertation is to explore the bridging function  $h$  that transforms the output embeddings from the visual encoder into a form compatible with the text decoder. Specifically, our goal is to find  $h$  such that the transformed visual embeddings  $Z_v = h(Y_v)$  can be seamlessly integrated into the text generation process of the decoder.

The simplest approach would be to directly use the output of the visual encoder  $Y_v$  as input to the text decoder, making  $h$  the identity function. This approach would only be feasible if the dimensions of the visual embeddings  $d_v$  and the text embeddings  $d_t$  are identical (i.e.,  $d_v = d_t$ ). However, even in such a case, the embeddings from the visual encoder would reside in a different latent space from the text embeddings, meaning that the text decoder would not properly interpret the visual embeddings.

Therefore, it is crucial to design a bridging function  $h$  that can map visual embeddings  $Y_v$  into the same semantic space as the text embeddings  $Z_t$ . This will allow the text decoder  $g$  to process the concatenated sequence  $Z = [Z_v, Z_t]$  and generate the final output text  $T_{output}$ .

In mathematical terms, the complete process can be described as:

$$g(Z = [Z_t, h(Y_v = f(I))]) = T_{output}, \quad (1.4)$$

where  $f$  is the visual encoder,  $g$  is the text decoder (both typically frozen), and  $h$  is the bridging function that is trained to enable smooth modality interaction.

This dissertation focuses on investigating different architectures and strategies for the bridging function  $h$ , with variations in the number of visual embeddings  $N$  and the dimensionality  $M_v$ , in order to optimize the performance of the vision-and-language model.

### 1.3 Contributions

The main contributions of this dissertation are:

- A comprehensive literature review detailing various approaches to bridging vision models with natural language models in Vision and Language systems.
- The development of European Portuguese Vision and Language datasets through machine translation of widely used English datasets.
- The creation of a new European Portuguese Vision and Language model designed for multimodal tasks.
- Submission of an article presenting the European Portuguese Vision and Language datasets and model.
- An ablation study on different bridging functions for Vision and Language models with frozen encoders and decoders, evaluating their impact on model performance.

## 1.4 Document Structure

- **Chapter 2, Background:** In this chapter we introduce some of the base concepts important to understanding the rest of the document. This includes the base elements used to process images and text and how these are trained.
- **Chapter 3, Related Work:** In this chapter we present some of the most important Vision and Language architectures, exploring the different possibilities for encoder models and the different decoder state-of-the-art models.
- **Chapter 4, V&L Tasks and Datasets:** In this chapter we present some of the different datasets and tasks used to train and evaluate Vision and Language models.
- **Chapter 5, A Portuguese Vision and Language Model:** In this chapter we present the paper we submitted, presenting a new European Portuguese Generative Vision and Language model and new Portuguese datasets by translating commonly used English datasets using machine translation.
- **Chapter 6, Ablation Study: CLS based models:** In this chapter we present an ablation study following the FROMAGe [14] architecture where we extract CLS tokens from the visual encoder and study how the model behaves by changing the layer from we extract, in how many tokens we transform it and how we can use CLS tokens from different layers together.
- **Chapter 7, A Ablation Study: All tokens Based Models:** In this chapter we present an ablation study similar to the previous one but where instead of using only CLS tokens we use all tokens from a given layer. We try different strategies to transform and reduce these tokens, ending up with a cross-attention mechanism from which we then explore the number of tokens and ViT layer.

- **Chapter 8, Conclusions:** In this chapter we present our conclusions with our main takeaways and possible future work.

## BACKGROUND

### 2.1 Elements of Vision and Language Model Architectures

Before exploring how multimodal models handle image and text processing together, it is essential to first understand how each modality is processed individually. This section focuses on the core architectures used for single-modality processing, the Transformer and Convolutional Neural Networks.

#### 2.1.1 Transformer

The most common algorithm used to encode text is the transformer, in this section we will explain step by step how it works.

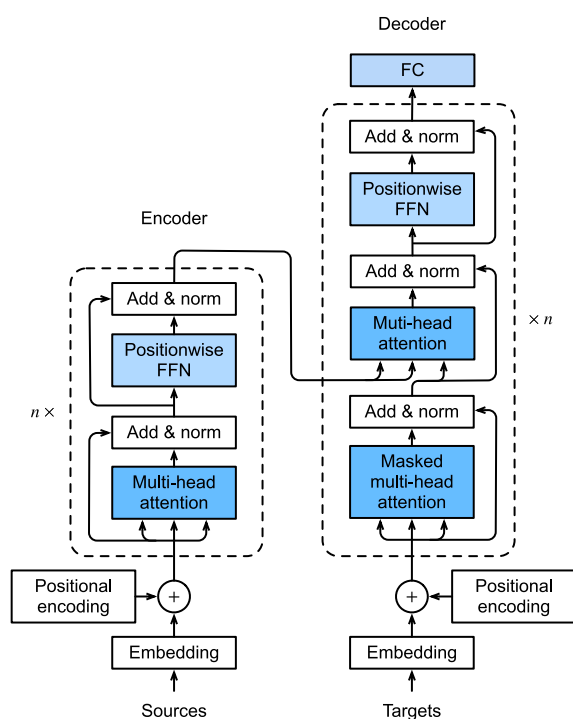


Figure 2.1: Architecture of the Original Transformer, with an encoder and a decoder. [43].

The original Transformer [37] is composed of two parts, an encoder which will transform the text into embeddings, and a decoder which will create new text as we can see in figure 2.1. Most recently, it is common to use only a decoder or only an encoder depending on the use.

### 2.1.1.1 Tokens and Embeddings

The first problem we have when trying to process text is its unstructured nature, in order to give it a structure the text needs to be tokenized, this means that the text will be broken into words, terms, sentences, symbols, or some other meaningful elements called tokens. Given a text  $T$ , a tokenizer will transform it into a sequence of these so called tokens  $(t_0, t_1, \dots, t_N)$ .

Now in order for a model to fully understand these tokens, they will be transformed into embeddings which are vector representations that capture semantic information about the tokens. These vectors are designed to encode relationships in a way that a machine learning model can understand and learn from.

After replacing the tokens with the corresponding embeddings we have a vector sequence  $X = [x_0, x_1, \dots, x_N] \in \mathbb{R}^{N \times d}$  where  $N$  corresponds to the number of tokens the text has been broken into and  $d$  the dimension of the vectors.

### 2.1.1.2 Positional embeddings

In language, the order and position of words are important to understand a sentence. In order to give the position information of each token in a sequence, positional embeddings are used. These will be learned vectors  $X^{pos} = [x_0^{pos}, x_1^{pos}, \dots, x_N^{pos}] \in \mathbb{R}^{N \times d}$ , each vector  $x_n^{pos}$  representing the n-th position.

Now we add these positional embeddings to our initial sequence of vector  $X$ , obtaining the following sequence of vector which will be the input of the transformer:

$$Y^0 = X + X^{pos} = [x_0 + x_0^{pos}, x_1 + x_1^{pos}, \dots, x_N + x_N^{pos}] \quad (2.1)$$

### 2.1.1.3 Self-attention

The most important feature of the transformer is self-attention, it allows the model to compute each token with respect to the others. Let  $Q^l = XW^{Ql}$ ,  $K^l = XW^{Kl}$  and  $V^l = XW^{Vl}$ , where  $W^{Ql} \in \mathbb{R}^{d \times d_q}$ ,  $W^{Kl} \in \mathbb{R}^{d \times d_q}$  and  $W^{Vl} \in \mathbb{R}^{d \times d_v}$  are trainable weight matrices for layer  $l$ , we have  $Q^l \in \mathbb{R}^{N \times d_q}$ ,  $K^l \in \mathbb{R}^{N \times d_q}$  and  $V^l \in \mathbb{R}^{N \times d_v}$ . In order to avoid a too confusing notation, we will be using these matrices without the  $l$  but it is important to note that this will work inside some layer. We define Attention as:

$$Att(Q, K, V) = \omega(QK^T) V \quad (2.2)$$

where  $\omega(\bullet) = softmax(\bullet / \sqrt{d_q})$ . This means that if we have  $O = Att(Q, K, V)$ , let  $q_j, k_j, v_j, o_j$  be j-th rows from  $Q, K, V$ , and  $O$  respectively then  $o_j$  will be a weighted average of

all the  $v_i$ 's where the weight of each  $v_i$  will be decided by the inner product of  $q_j$  with  $k_i$ . In practice, this means that for each token  $x_j$  we will output a token  $o_j$  that took attention from all the tokens.

#### 2.1.1.4 Multi Head Attention

Since there are multiple relationships that the tokens can have between each other, a multi-head attention architecture is used to concatenate the output of multiple self-attentions. In equation 2.3 we see how this is computed using  $W^O \in \mathbb{R}^{Hd_v \times d}$  to reduce the dimensionality back to  $d$ ,  $\oplus$  is denoting concatenation and  $h$  identifies the head.

$$MHA(X) = (O_1 \oplus O_2 \oplus \dots \oplus O_H)W^O \quad (2.3a)$$

$$O_h = Att(Q, K, V), \quad (2.3b)$$

where  $H$  is the total number of heads.

#### 2.1.1.5 Normalization Layer

Each transformer layer will be composed of a multi-head attention layer and a feedforward network layer, each of these layers will be followed by layer normalization with a residual layer. A residual layer consists of adding the output of a layer to the input, for an input  $X$  and a function  $f$  we would have  $f(X) + X$ .

Layer Normalization consists of normalizing the output of a layer getting  $LN(X) = \frac{X - \mu}{\sigma}$  where  $\mu$  and  $\sigma$  correspond, respectively, to the mean and the standard deviation. This way, for each Transformer layer we will have two trainable matrices  $W_1$  and  $W_2$ , obtaining the final output of the transformer layer as

$$M = MAB(X) = LN(X + MHA(X)) \quad (2.4a)$$

$$FFB(M) = LN(M + ReLU(MW_1)W_2) \quad (2.4b)$$

This way we get that  $Y^{l+1} = FFB(MAB(Y^l))$ .

#### 2.1.1.6 Contextual embeddings

As explained earlier, text embedding involves transforming words into vectors that represent their meanings in relation to other words. Some algorithms only perform this initial step, using a dictionary-like approach where each word is transformed into a static vector, independent of the surrounding context. While these algorithms have proven useful, they do not account for contextual information, as they assign a single representation to each word, regardless of how it's used in different contexts. However, it is crucial to acknowledge that a word can have multiple meanings. For instance, the word **left** in the sentence "*I left my phone on the left side of the table.*" carries two distinct meanings. This

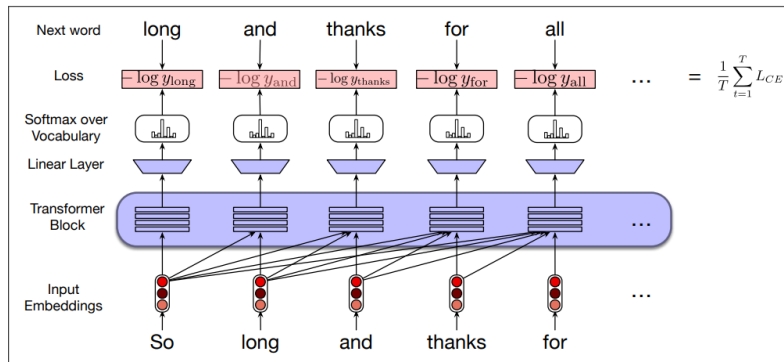


Figure 2.2: Structure of a transformer decoder. [11]

is where the transformer architecture becomes essential. By using mechanisms such as **self-attention** and **positional embeddings**, transformers can capture the context in which a word appears, producing what we call **contextual embeddings**. These embeddings provide a more nuanced representation of words based on their specific usage within a sentence.

### 2.1.1.7 Decoder

Until here we described how the transformer encoder works, now we will see how the decoder works. A transformer decoder works similarly to the transformer encoder but with some important differences. First of all, instead of self-attention, we have causal self-attention meaning that each token will pay attention to itself and the previous tokens but not to the preceding tokens. This happens because it predicts the next token with the present tokens, using softmax over the vocabulary to get the probabilities of each token, choosing the one with the highest probability. The selected token is then attached to the input to generate the next token until it has generated the answer, as seen in figure 2.2.

Instead of predicting always the token with the highest probability many decoders use beam search which takes the top  $n$  tokens to compute the next top  $n$  tokens until some number of tokens ahead, which will be the beam size, and then choose the token with which the combination with the highest probability begins.

## 2.1.2 Image Encoders

Next, we aim to apply a similar process to images, transforming them into vectors that can effectively represent the content within the image. It is important to note that an image can be seen as a tensor  $I \in \mathbb{R}^{H \times W \times C}$  where  $H$  is the height and  $W$  the width, each pair  $(h, w)$  corresponds to a pixel usually with 3 channels ( $C = 3$ ) using to the RGB system. The following methods will work with this tensor.

### 2.1.2.1 Convolutional Neural Networks

One of the most important and most used algorithms for Image processing is the Convolutional Neural Network (CNN) [16]. CNNs are very useful because their structure has inductive biases for images by introducing components that take into account the structure of images.

The main component of CNNs is the convolutional layer, implementing convolution operations which are composed of kernels that work as sliding windows that will multiply point-wise and sum the values of the window, as shown in figure 2.3.

Input		Kernel		Output																	
<table border="1" style="border-collapse: collapse; width: 60px; height: 60px;"> <tr><td style="padding: 2px 10px;">0</td><td style="padding: 2px 10px;">1</td><td style="padding: 2px 10px;">2</td></tr> <tr><td style="padding: 2px 10px;">3</td><td style="padding: 2px 10px;">4</td><td style="padding: 2px 10px;">5</td></tr> <tr><td style="padding: 2px 10px;">6</td><td style="padding: 2px 10px;">7</td><td style="padding: 2px 10px;">8</td></tr> </table>	0	1	2	3	4	5	6	7	8	*	<table border="1" style="border-collapse: collapse; width: 60px; height: 60px;"> <tr><td style="padding: 2px 10px;">0</td><td style="padding: 2px 10px;">1</td></tr> <tr><td style="padding: 2px 10px;">2</td><td style="padding: 2px 10px;">3</td></tr> </table>	0	1	2	3	=	<table border="1" style="border-collapse: collapse; width: 60px; height: 60px;"> <tr><td style="padding: 2px 10px;">19</td><td style="padding: 2px 10px;">25</td></tr> <tr><td style="padding: 2px 10px;">37</td><td style="padding: 2px 10px;">43</td></tr> </table>	19	25	37	43
0	1	2																			
3	4	5																			
6	7	8																			
0	1																				
2	3																				
19	25																				
37	43																				

Figure 2.3: Computation of a convolution [43].

There are other components used in CNN architectures like Average pooling and max pooling that downsample the convolution operation by averaging the values in the window or getting the maximum of the values, reducing the output dimension. These possible components of CNNs work well with images because besides maintaining the image structure they have important properties like translation invariance, meaning that some object in the image will produce the same output wherever it is.

### 2.1.2.2 Linear encoding of Image Tiles

Another way of encoding images is with a transformer, using image tiles as tokens, this is called Vision Transformer (ViT) [10]. Let  $I \in \mathbb{R}^{H \times W \times C}$  be an image, where  $H$  is the height,  $W$  is the width and  $C$  is the number of channels. This image is transformed as a sequence of flattened tiles of the image  $I_{patches} \in \mathbb{R}^{N \times (P^2 C)}$ , where  $(P, P)$  is the resolution of each patch and  $N = HW/P^2$  the length of our sequence of tokens. These patches are processed by a linear transformation, concatenated with a special class ([CLS]) token, and summed positional embeddings to transform them in the inputs of our transformer,  $X \in \mathbb{N} + \# \times \approx$ . This is illustrated in figure 2.4, and as we can see the final embedding of this special token is usually used as a representation of the whole image, being used for tasks like classification as in the figure.

A model also used is a hybrid of those two, where features are extracted from a CNN and then the patches are created from the resulting image. Even though CNNs have the advantage of having an inductive image bias, when trained with a lot of data the Visual Transformers perform on par with CNNs but are much more efficient since most of the computations can be done in parallel, for this reason, most VQA SOTA models use ViT.

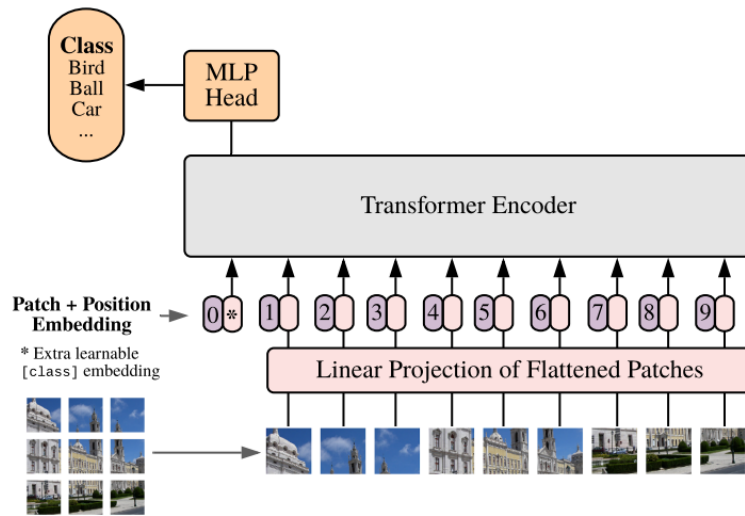


Figure 2.4: Visual Transformer Architecture [10]

## 2.2 Learning

Now we will see how to train Vision and Language models. We need loss functions that will be minimized, tasks that this minimization aims to achieve, and algorithms to minimize the loss. This will be done using batches created with data from some dataset. Therefore we need all these ingredients in order to properly train a model, we will now see how each of these works.

### 2.2.1 Loss functions

Before presenting the various loss functions we present a common notation in order to better understand the differences and similarities between the losses.

We will use  $X = \{x_1, \dots, x_N\}$  and  $Y = \{y_1, \dots, y_M\}$  from our data, this can represent observations and its labels/classes or different observations that we want compare. If  $x$  is of class  $y$  or  $x$  and  $y$  are similar we call  $(x, y)$  a positive pair, otherwise is a negative pair.

#### 2.2.1.1 Cross-Entropy Loss

Most models whose purpose is to output a token or a label to obtain these outputs, first obtain a vector where each value corresponds to a probability (between 0 and 1) of outputting a certain label/token, this means given an input  $x$  the model will give a probability  $p(y|x)$  for each possible answer  $y$ .

The cross-entropy loss compares the given probabilities with the ground truth, usually a vector full of 0's but a 1 in the place corresponding to the right answer.

Cross-entropy loss increases as the predicted probability diverges from the actual label. Given an input  $x$  and a label  $y$ , in the case where we only have two classes, this is  $y = 1$  or

$y = 0$  cross-entropy can be computed as:

$$H(x, y) = y \log(p_\theta(y|x)) + (1 - y) \log(1 - p_\theta(y|x)) \quad (2.5)$$

This loss allows us to train the parameters of this probability function  $p_\theta$  in a way that it is close to 0 if  $y$  is 0 and close to 1 if  $y$  is 1.

In the case where we have more than two classes cross-entropy will be computed as:

$$H(\phi, p) = - \sum_{y=1}^M \phi_{x,y} \log(p(y|x)) \quad (2.6)$$

where  $M$  is the number of classes,  $\phi$  the binary indicator (0 or 1) if the class label  $y$  is the correct classification of observation  $x$ . Since  $\phi$  is binary, being always 0 except when we  $y$  takes the right label  $y^+$  we can rewrite it as:

$$H(x) = - \log(p(y^+|x)) \quad (2.7)$$

### 2.2.1.2 Softmax loss

In order to compute the above probability  $p(y^+|x)$  we can use the softmax function as follows:

$$p(y^+|x) = \frac{e^{f_\theta(x, y^+)/\tau}}{\sum_{y \in Y} e^{f_\theta(x, y)/\tau}} \quad (2.8)$$

This way we train a function  $f_\theta : X \times Y \rightarrow \mathbb{R}$ , we can also represent this function as  $f_\theta : X \rightarrow \mathbb{R}^M$ , giving us for each label  $y$  a value that will be transformed this way in a probability. This function will take the largest values when the  $(x, y)$  is positive, i.e.  $x$  is of class  $y$ , and take small values otherwise. Using the exponential we grant that the values are positive and dividing by the summation we grant that the probabilities sum to 1. The hyper-parameter  $\tau$  is called temperature, for higher temperatures ( $\tau \rightarrow \infty$ ) this gives us a uniform probability distribution, for lower temperatures ( $\tau \rightarrow 0$ ) the probability of the class  $y$  corresponding to the higher value in  $f$  has a probability of 1. This way we end up with the following loss:

$$\mathcal{L}_{\text{softmax}}(x, y^+, Y, \theta) = - \log \frac{e^{f_\theta(x, y^+)/\tau}}{\sum_{y \in Y} e^{f_\theta(x, y)/\tau}} \quad (2.9)$$

where once again  $y^+$  represents the true label of  $x$ .

### 2.2.1.3 Constrastive loss

The goal of contrastive loss is for similar observations to be close and different observations to be far from each other in the embedding space. Here  $Y$  will be observations and not labels, in some cases  $X = Y$ . To do this we will train 2 functions  $g_\theta(\cdot) : X \rightarrow \mathbb{R}^d$  and  $h_\theta(\cdot) : Y \rightarrow \mathbb{R}^d$ , that will put  $X$  and  $Y$  in the same embedding space (in the case that  $X = Y$ ,  $g = h$ ). We want these functions to be in a way such that  $\|f(x_i) - g(y_j)\| \approx 0$  if

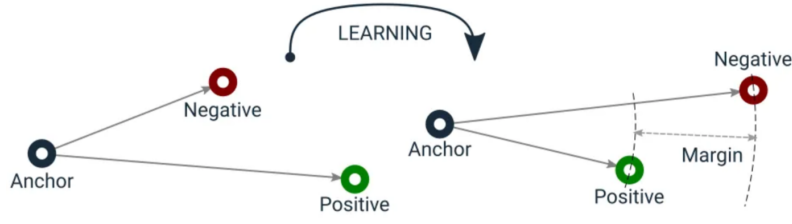


Figure 2.5: Illustration of the triplet loss. By minimizing the loss, the negative pair stays farther than the positive one from the anchor by the margin.

$(x_i, y_j)$  is a positive pair and  $\|f(x_i) - g(x_j)\| \geq \varepsilon$  if  $(x_i, y_j)$  is a negative pair, where  $\varepsilon > 0$  is a hyper-parameter. This way we get the following loss function:

$$\mathcal{L}_{\text{cont}}(x_i, y_j, \theta) = \phi_{x_i, y_j} d(g_\theta(x_i), h_\theta(y_j)) + \mathbb{1}(1 - \phi_{x_i, y_j}) \max(0, \varepsilon - d(g_\theta(x_i), h_\theta(y_j))) \quad (2.10)$$

where once again  $\phi_{x_i, y_j}$  is a binary function representing whether  $(x_i, y_j)$  is a positive pair.

#### 2.2.1.4 Triplet Loss

The triplet loss is one of the most popular loss functions for supervised similarity. It is similar to the contrastive loss as it also encourages dissimilar objects to be distant from each other and similar objects to be close. This time our margin  $\varepsilon$  represents the minimum difference we want between similar and dissimilar objects. With this idea in mind, this loss is used with triplets  $(x, y^+, y^-)$  where we call  $x$  the anchor ( $x, y^+$ ) is a positive pair and  $(x, y^-)$  is a negative pair, as we can see in figure 2.5. As we want  $y^+$  to be closer to  $x$  than  $y^-$  we will have the following loss:

$$\mathcal{L}_{\text{triplet}}(x, y^+, y^-, \theta) = \max(0, d(g_\theta(x), h_\theta(y^+)) - d(f_\theta(x), f_\theta(y^-)) + \varepsilon) \quad (2.11)$$

#### 2.2.1.5 Lifted Structured Loss

Lifted Structured Loss follows the triplet loss idea but uses all positive and negative pairs at a time instead of a triplet with an anchor a positive and a negative. Let  $\mathcal{P}$  and  $\mathcal{N}$  be the set of positive pairs and the set of negative pairs respectively we want the distance between positive pairs to be close to 0 and the distance between negative pairs to be at least  $\varepsilon$ . Let  $D_{ij} = d(g_\theta(y_i), h_\theta(y_j))$ , for each positive pair  $(i, j)$  we want to minimize  $D_{ij}$  and maximize  $D_{i,k}$  or  $D_{j,k}$  where  $k$  is the negative observation closest to either  $i$  or  $j$ . This way we obtain the following loss:

$$\mathcal{L}_{\text{struct}} = \frac{1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \max(0, \mathcal{L}_{\text{struct}}^{(ij)})^2 \quad (2.12)$$

where  $\mathcal{L}_{\text{struct}}^{(ij)} = D_{ij} + \max\left(\max_{(i,k) \in \mathcal{N}} \varepsilon - D_{ik}, \max_{(j,l) \in \mathcal{N}} \varepsilon - D_{jl}\right)$

The red part is used for mining hard negatives. However, it is not smooth and may cause the convergence to a bad local optimum in practice. Thus, it is relaxed to be:

$$\mathcal{L}_{\text{struct}}^{(ij)} = D_{ij} + \log \left( \sum_{(i,k) \in \mathcal{N}} \exp(\varepsilon - D_{ik}) + \sum_{(j,l) \in \mathcal{N}} \exp(\varepsilon - D_{jl}) \right) \quad (2.13)$$

### 2.2.1.6 Supervised Contrastive loss

Supervised contrastive loss generalizes the triplet loss to include comparisons between multiple negative pairs. This way instead of an anchor, one positive observation and one negative observation we use an anchor  $x$ , one positive  $y^+$  and  $N - 1$  negatives  $Y^- = \{y_1^-, \dots, y_{N-1}^-\}$ . We get the following loss:

$$\mathcal{L}_{\text{N-pair}}(x, x^+, X^-) = -\log \frac{\exp(g(x)^\top h(y^+)/\tau)}{\exp(g(x)^\top h(y^+)/\tau) + \sum_{i=1}^{N-1} \exp(g(x)^\top h(y_i^-)/\tau)} \quad (2.14)$$

we can now observe that if take  $f_\theta(x, y) = g(x)^\top h(y)$  which computes a similarity score between  $x$  and  $y$  we get the softmax loss. The difference between these two being that for the softmax  $y$  represents a label and here  $y$  represents an observation.

### 2.2.1.7 NCE

Noise Contrastive Estimation (NCE) uses logistic regression in order to separate the target answer from noise, which correspond to the positive and negative observations. Given some context  $x$ , we have  $y^+ \sim p_\theta(y^+|x)$ . For the negative observation, we have  $y^- \sim q(y^-)$  where  $q$  represents a noise distribution. Applying the logit function used by the logistic regression we obtain:

$$\ell_\theta(u) = \log \frac{p_\theta(u)}{q(u)} = \log p_\theta(u) - \log q(u) \quad (2.15)$$

Now we convert the logits into probabilities with the sigmoid function  $\sigma(\cdot)$  and apply the cross-entropy loss, obtaining the final loss:

$$\mathcal{L}_{\text{NCE}} = -\frac{1}{N} \sum_{i=1}^N [\log \sigma(\ell_\theta(y_i^+)) + \log(1 - \sigma(\ell_\theta(y_i^-)))] \quad (2.16)$$

where  $\sigma(\ell) = \frac{1}{1 + \exp(-\ell)} = \frac{p_\theta}{p_\theta + q}$

### 2.2.1.8 InfoNCE

The InfoNCE loss is inspired by the NCE loss and uses categorical cross-entropy loss to identify the positive observation amongst a set of noise observations.

The main idea of this loss is that given the context  $x$  the positive observation should be drawn from  $p(y|x)$  and the negative ones from  $p(y)$ . This way by taking a sample  $Y = \{y_1, \dots, y_M\}$  with only one positive sample  $y^+$ , the probability of detecting the

positive sample correctly corresponds to the probability of observing this sample when  $y^+$  is drawn from  $p(y|x)$  normalized by the sum of the probabilities of each  $y_i$  being the one to be drawn from  $p(y|x)$ , this is:

$$p(C = \text{pos}|Y, x) = \frac{p(y^+|x) \prod_{i=1, \dots, M; y_i \neq y^+} p(y_i)}{\sum_{j=1}^M [p(y_j|x) \prod_{i=1, \dots, M; i \neq j} p(x_i)]} = \frac{\frac{p(y^+|x)}{p(y^+)}}{\sum_{j=1}^M \frac{p(y_j|x)}{p(y_j)}} = \frac{f(x, y^+)}{\sum_{j=1}^M f(x, y_j)} \quad (2.17)$$

where we train the function  $f(x, y) \propto \frac{p(x|c)}{p(x)}$ . Obtaining the following loss:

$$\mathcal{L}_{\text{InfoNCE}} = -E \left[ \log \frac{f(x_{\text{pos}}, c)}{\sum_{x \in X} f(x, c)} \right] \quad (2.18)$$

We can see the function  $f$  trained for the InfoNCE loss as a similarity function, if we take it as  $f(x, c) = e^{g(x)^\top h(y)/\tau}$  we get the Supervised Contrastive loss. This way although they seem to have different inspirations, we can see the InfoNCE as a generalization of the supervised contrastive loss.

## 2.2.2 Tasks

### 2.2.2.1 Masked Language Modeling (MLM)

Masked Language Modeling is one of the most common task for language modeling pre training but is also used for V & L models.

Let  $m \in \mathbb{N}^m$  where  $m$  is the number of masked tokens and  $m$  is the set of masked indices. Given an image-text pair some of the words will be randomly masked and the masked ones  $\tilde{w}_m$  will be replaced with a special token [MASK].

The goal is to train the model to predict the masked words based on the other words  $\tilde{w}_{\setminus m}$  and the paired image  $\tilde{v}$  by minimizing the loss corresponding to the negative log-likelihood:

$$\mathcal{L}_{MLM}(\theta) = -\mathbb{E} \left[ \log P_\theta(\tilde{w}_m | \tilde{w}_{\setminus m}, \tilde{v}) \right] \quad (2.19)$$

### 2.2.2.2 Image Text Matching (ITM)

In ITM, given a batch of matched or mismatched image-caption pairs, the model needs to identify which images and captions correspond to each other. This way, given an image-caption pair  $(\tilde{v}, \tilde{w})$  the model will output a score  $s_\theta(\tilde{v}, \tilde{w})$ , here it is used the cross-entropy loss for two classes (1 if the pair matches, 0 if it doesn't), obtaining the following:

$$\mathcal{L}_{ITM}(\theta) = -\mathbb{E} \left[ y \log(s_\theta(\tilde{v}, \tilde{w})) + (1 - y) \log(1 - s_\theta(\tilde{v}, \tilde{w})) \right] \quad (2.20)$$

### 2.2.2.3 Contrastive Learning

Contrastive learning is a commonly used pre training task that aims to map the image and the text into the same feature map, this means that a picture of a dog and the word "dog" should be represent by similar vectors. This usually works using an image encoder

and a text encoder that will be trained from large datasets with image-caption pairs. A contrastive loss function is used to maximize the similarity of the embeddings of image-text pairs if they match and minimize the similarity if they don't. An example of a model that uses this training is CLIP [28] that uses a cosine similarity and we can see how this contrastive learning works in figure 3.4 where the values in blue are to be maximized and the ones in white to be minimized.

This way let  $v_i$  and  $w_j$  be the embeddings corresponding to image  $i$  and caption  $j$  we want these embeddings to be similar when  $i = j$  and not similar when  $i \neq j$ , let  $sim$  be a similarity function,  $s_{i,j}^{i2t} = sim(v_i, w_j)$  and  $s_{i,j}^{t2i} = sim(w_i, v_j)$  we want to minimize the following losses:

$$\mathcal{L}_{ITC}^{i2t}(\theta) = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(s_{i,i}^{i2t}/\sigma)}{\sum_{j=1}^N \exp(s_{i,j}^{i2t}/\sigma)} \quad (2.21)$$

$$\mathcal{L}_{ITC}^{t2i}(\theta) = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(s_{i,i}^{t2i}/\sigma)}{\sum_{j=1}^N \exp(s_{i,j}^{t2i}/\sigma)} \quad (2.22)$$

Now a common way to do it is to get the mean of these two losses:

$$\mathcal{L}_{ITC}(\theta) = \frac{1}{2} [\mathcal{L}_{ITC}^{i2t}(\theta) + \mathcal{L}_{ITC}^{t2i}(\theta)] \quad (2.23)$$

It is important to note that these losses correspond to the Supervised Contrastive Loss or the InfoNCE where equation 2.2.2.3 uses the images as the  $X$  and the text as the  $Y$  and equation 2.2.2.3 uses  $X$  as the image and  $Y$  as the text.

#### 2.2.2.4 Causal Language Generation

Causal Language Generation is the main task to train language models because it focuses on predicting the next word as we have seen previously about Large Language Models, looking back at figure 2.2 we see how this task is handled, by predicting the words one by one, taking only attention to the previous words. Given a sequence of words  $w = w_1, \dots, w_n$  the general training objective is to minimize the loss:

$$\mathcal{L}_{\text{causal}}(\theta) = \sum_{i=1}^n \log p_{\theta}(w_i | w_{<i}) \quad (2.24)$$

For vision and language models we can add some context  $c$  which can be some previous text, an image or both, obtaining:

$$\mathcal{L}_{\text{causal}}(\theta) = \sum_{i=1}^n \log p_{\theta}(w_i | w_{<i}, c) \quad (2.25)$$

### 2.2.3 Model Optimization

In order to minimize the loss function we need an optimizer, that is, an algorithm that tries to find the parameters that give us the lowest loss values. We will now see how these algorithms work in order to train our models.

### 2.2.3.1 Gradient Descent

Most optimizers are based on Gradient Descent which uses the gradient of the loss function, that is, the derivatives with respect to all the parameters. Let  $f$  be a function and  $\theta = [\theta_1, \dots, \theta_N]$  the inputs of our function. We want to find a  $\theta$  that minimizes  $f(\theta)$ . Using the first-order Taylor expansion we have that:

$$f(\theta + \epsilon) = f(\theta) + \epsilon^\top \nabla f(\theta) + \mathcal{O}(\|\epsilon\|^2). \quad (2.26)$$

Now we take  $\epsilon = -\eta \nabla f(\theta)$ :

$$f(\theta - \eta \nabla f(\theta)) = f(\theta) - \eta \nabla f(\theta)^\top \nabla f(\theta) + \mathcal{O}(\eta^2 \|\nabla f(\theta)\|^2). \quad (2.27)$$

with this idea in mind, we can predict that for an  $\epsilon$  small enough we have:

$$f(\theta - \eta \nabla f(\theta)) < f(\theta). \quad (2.28)$$

For gradient descent, we first choose an initial value for  $\theta$  (usually random) and a value for  $\eta$ , and each iteration of our algorithm we will have:

$$\theta \leftarrow \theta - \eta \nabla f(\theta). \quad (2.29)$$

### 2.2.3.2 Mini-batch Stochastic Gradient Descent

In order to compute the gradient we would need to compute it for every observation on our dataset. For a large dataset that would be very expensive and inefficient therefore we will take a sample from our data and obtain an estimation of the gradient from that sample. These samples correspond to the batches explained above. This way we obtain:

$$g_t = \partial_\theta \frac{1}{|\mathcal{B}_t|} \sum_{i \in \mathcal{B}_t} f(x_i, \theta) \quad (2.30)$$

where  $\mathcal{B}_t$  is our batch and  $x_i \in \mathcal{B}_t$  the observations from our batch. This way for each iteration we have:

$$\theta \leftarrow \theta - \eta_t g_t \quad (2.31)$$

where we have  $\eta_t$  instead of a constant  $\eta$  because we want our training rate to decrease with time in order to converge faster.

To this algorithm, we call Mini-Batch Stochastic Gradient Descent. When  $|\mathcal{B}_t| = 1$ , that is, the gradient is calculated only with one observation, we call it just Stochastic Gradient Descent.

### 2.2.3.3 Stochastic Gradient Descent with Momentum

To handle the variation introduced by our sampling we can add momentum, this way we introduce a velocity vector  $v$  that corresponds to a leaky average of the past gradients. We define it as:

$$v_t = \beta v_{t-1} + g_t \quad (2.32)$$

with  $\beta \in (0, 1)$ , this way we have:

$$v_t = \sum_{\tau=0}^{t-1} \beta^\tau g_{t-\tau} \quad (2.33)$$

For each iteration of our algorithm we will have:

$$\begin{aligned} v_t &\leftarrow \beta v_{t-1} + g_t, \\ \theta_t &\leftarrow \theta_{t-1} - \eta_t v_t. \end{aligned} \quad (2.34)$$

#### 2.2.3.4 Adagrad

This optimizer changes the learning rate taking into account the magnitude of the gradients, i.e., with larger gradients, we will have smaller learning rates. With this idea, we will accumulate the squares of the gradients into a variable  $s_t$  and we will define the learning rate as  $\frac{\eta}{\sqrt{s_t + \epsilon}}$ . Obtaining at each iteration:

$$\begin{aligned} s_t &= s_{t-1} + g_t^2, \\ \theta_t &= \theta_{t-1} - \frac{\eta}{\sqrt{s_t + \epsilon}} g_t. \end{aligned} \quad (2.35)$$

#### 2.2.3.5 RMSProp

Accumulating the squares of the gradients as in Adagrad makes  $s_t$  grow without bound, RMSProp changes slightly the way we compute  $s_t$  using a new parameter  $\gamma > 0$  obtaining:

$$\begin{aligned} s_t &= \gamma s_{t-1} + (1 - \gamma) g_t^2, \\ \theta_t &= \theta_{t-1} - \frac{\eta}{\sqrt{s_t + \epsilon}} g_t. \end{aligned} \quad (2.36)$$

#### 2.2.3.6 Adam

The Adam optimizer is one of the most used optimizers and fuses the ideas from the momentum as well as from the RMSProp, using both  $v_t$  and  $s_t$  defined as:

$$\begin{aligned} v_t &= \beta_1 v_{t-1} + (1 - \beta_1) g_t, \\ s_t &= \beta_2 s_{t-1} + (1 - \beta_2) g_t^2. \end{aligned} \quad (2.37)$$

with  $\beta_1, \beta_2 > 0$ . Common choices are  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . Now these vectors are normalized obtaining:

$$\hat{v}_t = \frac{v_t}{1 - \beta_1^t} \text{ and } \hat{s}_t = \frac{s_t}{1 - \beta_2^t}. \quad (2.38)$$

and finally we get  $\theta_t \leftarrow \theta_{t-1} - g'_t$  where

$$g'_t = \frac{\eta \hat{v}_t}{\sqrt{\hat{s}_t + \epsilon}}. \quad (2.39)$$

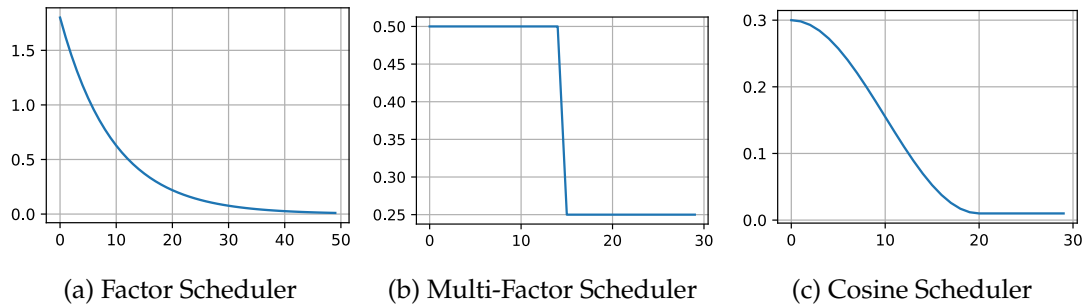


Figure 2.6: Comparison and examples of different scheduling functions. [43]

### 2.2.3.7 AdamW

We can add a weight decay in our optimizer in order to avoid overfitting. AdamW subtracts from our weights the term  $\lambda\theta$  in order to avoid large weight values. This way we maintain everything in the Adam optimizer except that:

$$\theta_t \leftarrow \theta_{t-1} \eta_t \left( \frac{\alpha \hat{v}_t}{\sqrt{\hat{s}_t + \epsilon}} + \lambda \theta_{t-1} \right) \quad (2.40)$$

where  $\eta_t$  is a scheduled multiplier.

### 2.2.3.8 Learning Rate Scheduling

In AdamW we used scheduled the learning rate  $\eta_t$ . This means that we use a predetermined function that gives us the learning rate at each timestep, usually at each epoch. As mentioned earlier to get a good convergence we want the learning rate at the end to be small, creating decreasing functions. Some of the most used functions are:

- Factor Scheduler -  $\eta_{t+1} \leftarrow \eta_t \cdot \alpha$  for  $\alpha \in (0, 1)$  or  $\eta_{t+1} \leftarrow \max(\eta_{\min}, \eta_t \cdot \alpha)$  to prevent the learning rate to be unreasonably small.
- Multi-Factor Scheduler - it works similar to Factor scheduling but in "steps" this means that we will have a set of timesteps and the decrease  $\eta_{t+1} \leftarrow \eta_t \cdot \alpha$  only happens when  $t \in s$ .
- Cosine Scheduler -  $\eta_t = \eta_T + \frac{\eta_0 - \eta_T}{2} (1 + \cos(\pi t/T))$  for  $t < T$  and  $\eta_t = \eta_T$  for  $t \geq T$ , this function decreases the learning rate slowly at the beginning, smoothly transitioning from the maximum  $\eta_0$  to the minimum  $\eta_T$ .

In figure 2.6 we can see how the mentioned functions behave. All these functions are based on the idea that we would start with larger learning rates that would decrease, but large initial learning rates may be a problem since expect significant divergence due to the random initialization of neural network weights. To address this we can add a warmup period in which the learning rate increases to its maximum, usually increasing at a linear

rate, only after this warmup period is the decreasing scheduling function applied. In figure 2.7 we can see how the cosine scheduler function would be with warmup.

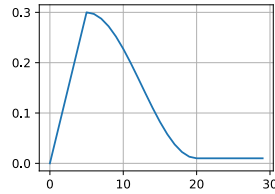


Figure 2.7: Example of a cosine scheduler function with a linear warmup. [43]

### 2.2.4 LoRA

Pre-training a model is computationally expensive and can take a long time. Therefore after the pre-training, we want fine-tuning not to be that expensive, but a full fine-tuning requires to that the derivative with respect to every parameter. For large models with billions of parameters this would still be very inefficient. For a weight matrix  $W_0 \in \mathbb{R}^{d \times r}$  we would get the derivatives with respect of each  $dr$  parameter and update the matrix obtaining  $W_0 + \Delta W$ . It was found that during fine-tuning this matrix  $\Delta W$  is not full rank having low intrinsic dimension. With this knowledge Low-Rank Adaptation (LoRA) fine tunes decomposes  $\Delta W$  into  $BA$  where  $B \in \mathbb{R}^{d \times r}$  and  $A \in \mathbb{R}^{r \times k}$  where  $r \ll \min(d, k)$ . This way we froze the matrix  $W_0$  and only train  $r(k + d) \ll kd$  parameters, making this training more efficient.

### 2.2.5 Batch creation

We cannot use all the observations of our dataset for each iteration of our training therefore we need to create batches. The most basic batch creation strategy is to choose a batch size  $N$  and choose  $N$  random observations every time, but this may have some problems for example:

1. Wrong format - depending on the task we may have to transform the data in order to have observations that can be processed by the loss function,
2. Bad sampling - sampling uniformly may be inefficient for the task at hand,
3. Not enough variation in the data - the dataset may be small or have little in-class variation requiring the use of data augmentation.

In tasks such as contrastive learning or when using the triplet loss we need to extract data from the dataset in such a way that we have an anchor, positive pairs and negative pairs.

### 2.2.5.1 Sampling

When sampling for a contrastive task we usually want for each positive pair many negative observations. When choosing the negatives uniformly we may have a lot of easy negatives, that is, negative observations that are far enough from the anchor, for the triplet loss this corresponds to  $d(g_\theta(x), h_\theta(y^+)) + \varepsilon < d(f_\theta(x), f_\theta(y^-))$ . We prefer to have more semi-hard negatives ( $d(g_\theta(x), h_\theta(y^+)) < d(f_\theta(x), f_\theta(y^-)) + \varepsilon$ ) or hard negatives ( $d(f_\theta(x), f_\theta(y^-)) < d(g_\theta(x), h_\theta(y^+))$ ) this way the training would be more efficient by taking more effort in separating the negatives from the positives.

This problem leads to hard negative mining and semi-hard negative mining, which corresponds to the creation of batches where we choose the negatives that are hard or semi-hard respectively. With these ideas, we can build many different sampling approaches. One problem regarding semi-hard and hard negative mining is that when the distance is closer to zero there is a higher effect of the noise increasing the variance of the gradient.

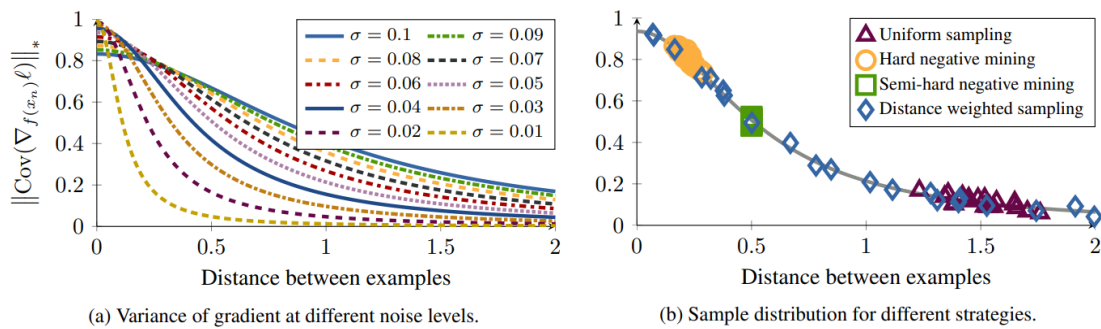


Figure 2.8: a) shows the nuclear norm of a noisy gradient estimate for various levels of noise (b) compares different sampling strategies. [41]

Figure 2.8 (a) shows this variance effect. Note that in these embedding approaches, the vectors are usually normalized, lying on the surface of a hypersphere, leading to the fact that by choosing two random points their expected distance is  $\sqrt{2}$ . [41] introduces a sampling strategy that aims to sample uniformly by distance by setting the probability of choosing some negative  $n$  given an anchor  $a$  as  $P(n^* = n|a) \propto \min(\lambda, q^{-1}(D_{an}))$  where we choose some  $\lambda$ ,  $q$  is the distribution of the distance between two random point and  $D_{an}$  the distance between  $a$  and  $n$ . This means that the probability of including  $n$  in the sample is proportional to the inverse of  $q$  but using  $\lambda$  to guarantee a minimum probability. Figure 2.8 (b) shows the difference in the distribution of distances between this sampling strategy and others used.

## 2.3 Summary

In this chapter, we presented an introduction to the components that are used to create Vision and Language models. We started by presenting how language is processed explaining in detail one of the most important algorithms in Natural Language Processing,

the Transformer. Then we explained one common algorithm to process images, the Convolutional Neural Network, and how many SOTA models also use the Transformer for image processing. After that, we focused on learning explaining some of the most important loss functions, tasks, and optimization algorithms as well as some important precautions to take when creating batches. All these different pieces are the base for creating and training good Vision and Language models.

## RELATED WORK

### 3.1 Vision and Language Model Encoder Architectures

We saw earlier how we can process text and image individually, now we will see different ways to process them together in a single Transformer. For this we will consider our image  $I$  and text  $T$  have already been pre-processed, transformed into embeddings, obtaining  $Y_{\mathcal{L}}^0 \in \mathbb{R}^{d \times N_l}$  and  $Y_{\mathcal{V}}^0 \in \mathbb{R}^{d \times N_v}$  as the language embeddings and vision embeddings respectively, where  $N_l$  and  $N_v$  correspond to the number of tokens of the text and the image respectively. This should correspond to the initial input embedding the text and the result of the linear encoding of image tiles.

These embeddings serve as the input to a Transformer, which is composed of  $L$  layers. Each layer in the Transformer processes the text and image embeddings and produces updated embeddings for both modalities. The output of layer  $l$  for text is represented as  $Y_{\mathcal{L}}^l$  and for images, it is  $Y_{\mathcal{V}}^l$ . This notation emphasizes that these outputs evolve as they pass through each Transformer layer. For simplicity, the embeddings are often written as  $Y_{\mathcal{L}}$  for text and  $Y_{\mathcal{V}}$  for images without explicitly indicating the layer. However, it's important to remember that these embeddings are updated at each layer within the Transformer.

#### 3.1.1 A Unified Mathematical Framework

In this section, we will define a mathematical framework based on [4] to better understand the different vision-language encoding architectures as different choices of variables in the same framework. We start by reshaping the way we use attention, by considering the use of both modalities. Instead of having the matrices  $W^Q$ ,  $W^K$ , and  $W^V$  we will have these three per modality, this way  $W_{h_{\mathcal{L}}}^Q$  is the query weight matrix for the text and  $W_{h_{\mathcal{V}}}^Q$  the query weight matrix for the image. With this in mind, we will set the

$$S_{AB} = Q_A K_B \tag{3.1}$$

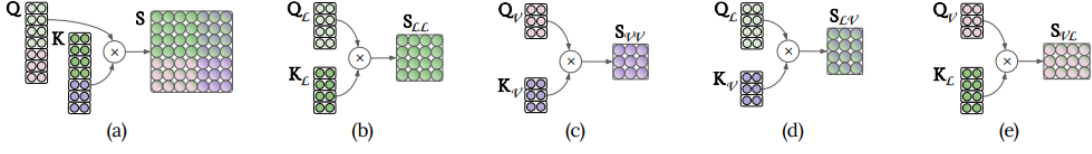


Figure 3.1: Visualization of the different score matrices [4]

where  $Q_A = Y_A W_A^Q$  and  $K_B = Y_B W_B^K$  are queries and the keys coming respectively from modality  $A$  and modality  $B$ . With this, we define

$$S_\gamma = \begin{pmatrix} \varepsilon^{\gamma_{LL}} S_{LL} & \varepsilon^{\gamma_{LV}} S_{LV} \\ \varepsilon^{\gamma_{VL}} S_{VL} & \varepsilon^{\gamma_{VV}} S_{VV} \end{pmatrix} \quad (3.2)$$

where  $\varepsilon = -\infty$ ,  $\gamma_{AB} \in \{0, 1\}$  therefore  $\varepsilon^{\gamma_{AB}} \in \{-\infty, 1\}$ . We have that when  $\gamma_{AB} = 0$ ,  $\varepsilon^{\gamma_{AB}} = 1$  preserving  $S_{AB}$  and when  $\gamma_{AB} = 1$ ,  $\varepsilon^{\gamma_{AB}} = -\infty$  setting the contribution of  $S_{AB}$  as 0 when using the softmax. Here we have  $A, B \in \{\mathcal{V}, \mathcal{L}\}$ , for a better understanding of  $S_{AB}$  see figure 3.1. With this we will have the binary variables  $\gamma = \{\gamma_{LV}, \gamma_{LV}\gamma_{VL}, \gamma_{LL}, \gamma_{VV}\} \in \{0, 1\}^4$  these values act as gates that regulate the cross-modal interactions within a layer. Now we are ready to define multimodal multi-headed attention as the linear projection of a concatenation of single self-attention blocks.

$$MHA(Y_{\mathcal{L}}, Y_{\mathcal{V}}) = [O_1 \oplus \dots \oplus O_H] \begin{pmatrix} W_{\mathcal{L}}^O \\ W_{\mathcal{V}}^O \end{pmatrix} \quad (3.3)$$

where

$$\begin{aligned} O_h &= Att \left( \left( \begin{pmatrix} Y_{\mathcal{L}} W_{h\mathcal{L}}^Q \\ Y_{\mathcal{V}} W_{h\mathcal{V}}^Q \end{pmatrix}, \begin{pmatrix} Y_{\mathcal{L}} W_{h\mathcal{L}}^K \\ Y_{\mathcal{V}} W_{h\mathcal{V}}^K \end{pmatrix}, \begin{pmatrix} Y_{\mathcal{L}} W_{h\mathcal{L}}^V \\ Y_{\mathcal{V}} W_{h\mathcal{V}}^V \end{pmatrix}; \gamma \right) \\ &= Att \left( \left( \begin{pmatrix} Q_{h\mathcal{L}} \\ Q_{h\mathcal{V}} \end{pmatrix}, \begin{pmatrix} K_{h\mathcal{L}} \\ K_{h\mathcal{V}} \end{pmatrix}, \begin{pmatrix} V_{h\mathcal{L}} \\ V_{h\mathcal{V}} \end{pmatrix}; \gamma \right) \\ &= \omega(S_{h\gamma}) \begin{pmatrix} V_{h\mathcal{L}} \\ V_{h\mathcal{V}} \end{pmatrix} \end{aligned} \quad (3.4)$$

With this formulation we can build the layers in figure 3.2 where (a), (b) and (c) are particular cases of layer (d) we just described. This way we can define an encoder layer as

$$(Y_{\mathcal{L}}^l, Y_{\mathcal{V}}^l) = f_l((Y_{\mathcal{L}}^{l-1}, Y_{\mathcal{V}}^{l-1}), \gamma^l, \tau^l) \quad (3.5)$$

where  $f_l$  represents a transformer block as formulated in Eq. 2.4 but following the particularities of the chosen  $\gamma^l$  and  $\tau^l$ .

### 3.1.2 Single-stream attention

Single stream architectures work by concatenating the visual and textual embeddings in one sequence of tokens and using it as input of the transformer, this can be seen as

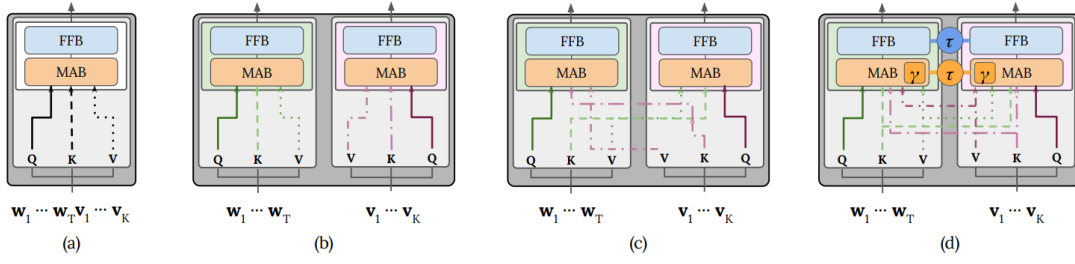


Figure 3.2: Encoder architectures [4]

setting  $\gamma_{\mathcal{L}\mathcal{V}} = \gamma_{\mathcal{V}\mathcal{L}} = \gamma_{\mathcal{L}\mathcal{L}} = \gamma_{\mathcal{V}\mathcal{V}} = 0$ . Corresponds to layer (a) in figure 3.2 using the score matrix in figure 3.1 (a). One example of this is the ViLT model [12] which also uses class embeddings  $x_{\mathcal{L}}^{class}$  and  $x_{\mathcal{V}}^{class}$  before each modality and type embeddings  $X_{\mathcal{L}}^{type}$  and  $X_{\mathcal{V}}^{type}$  to indicate to which modality belongs each token, as can be seen in figure 3.3 and in equation 3.6. Having The multi-head attention is then applied to the sequence  $z_0$ .

$$\bar{X}_{\mathcal{L}} = [x_{\mathcal{L}}^{class}; x_{\mathcal{L}}^1; \dots; x_{\mathcal{L}}^L] + X_{\mathcal{L}}^{pos} \quad (3.6a)$$

$$\bar{X}_{\mathcal{V}} = [x_{\mathcal{V}}^{class}; x_{\mathcal{V}}^1; \dots; x_{\mathcal{V}}^L] + X_{\mathcal{V}}^{pos} \quad (3.6b)$$

$$Y^0 = [\bar{X}_{\mathcal{L}} + X_{\mathcal{L}}^{type}; \bar{X}_{\mathcal{V}} + X_{\mathcal{V}}^{type}] \quad (3.6c)$$

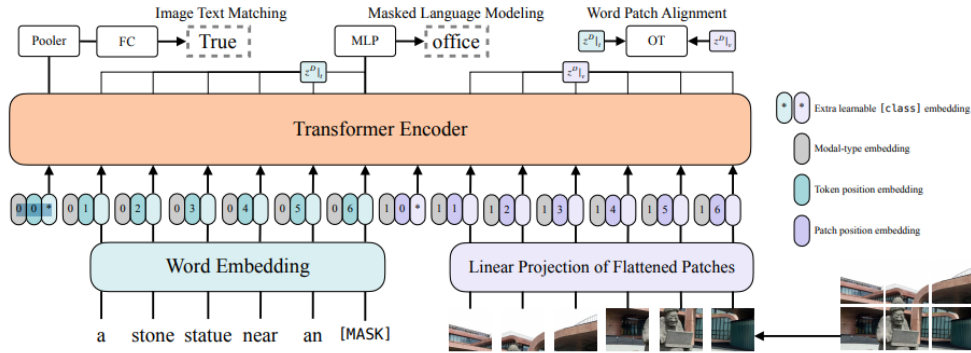


Figure 3.3: The single stream architecture of the ViLT model. [12]

As we can see ViLT is trained using Image Text Matching (IMT) and Masked Language Modeling (MLM) as explained in section 2.2.2 and with an extra Word Patch Alignment (WPA) task which computes the alignment between two subsets of the set of text tokens and the set of image tokens.

### 3.1.3 Dual-stream attention

Dual stream models encode the modalities separately using layers like (b) in figure 3.2 that corresponds to setting  $\gamma_{\mathcal{L}\mathcal{V}} = \gamma_{\mathcal{V}\mathcal{L}} = 1$  and  $\gamma_{\mathcal{L}\mathcal{L}} = \gamma_{\mathcal{V}\mathcal{V}} = 0$ , using the score matrices

in figures 3.1 (b) and 3.1 (c). One example of this is CLIP [28] which uses cosine similarity to connect the encodings, CLIP compares text with image and can be used for VQA with an image and a set of answers returning the most similar answer as can be seen in figure 3.4.

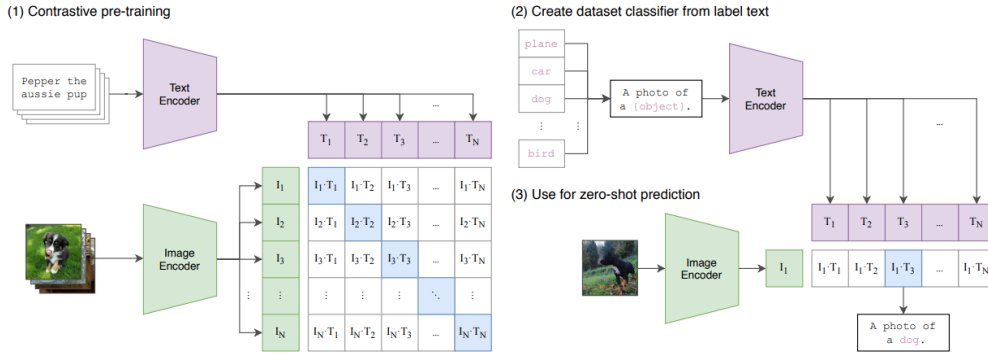


Figure 3.4: The dual stream architecture of the CLIP model. [28]

As we can see in figure 3.4 the pre-training of CLIP is very simple, using only the Contrastive Learning task as explained in section 2.2.2, using the similarity of the text and images CLS tokens.

### 3.1.4 Cross-attention

Another kind of layer architecture is the cross-attention, here we use the queries from one modality with the keys and values of the other modality, this is the layer (c) in figure 3.2 that corresponds to setting  $\gamma_{\mathcal{L}\mathcal{V}} = \gamma_{\mathcal{V}\mathcal{L}} = 0$  and  $\gamma_{\mathcal{L}\mathcal{L}} = \gamma_{\mathcal{V}\mathcal{V}} = 1$ , using the score matrices in figures 3.1 (d) and 3.1 (e). With this we will compute the new image embeddings using the Q from the image and the K and V from the text and the text embeddings will be computed using the Q from the text and the K and V from the image.

Many models use combinations of these layers, or others that can be made by changing  $\gamma$ . One example is LxMERT [32] which starts with a dual stream architecture and afterward a set of cross-attention layers intercalated with intramodal dual stream self-attention, as we see in figure 3.5.

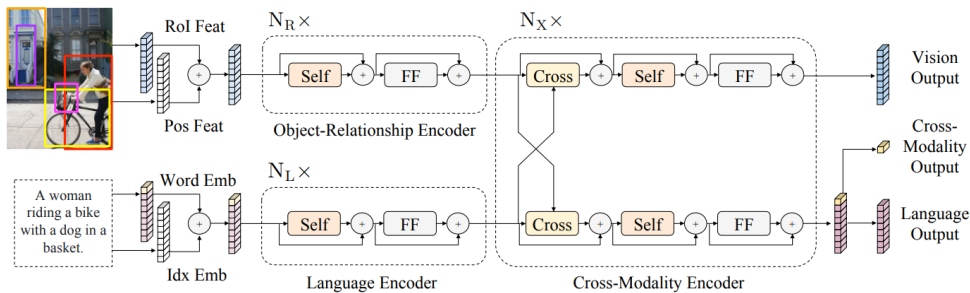


Figure 3.5: The cross-attention architecture of the LXMERT model. Note that the initial layers are dual-stream. [32]

LxMERT trains with the MLM task using the Language output we see in figure 3.5 and ITM using the Cross Modality output.

### 3.1.5 Applications of V&L Encoder Models

These encoder-only models while presenting good encoding capabilities are not able to generate text, therefore they cannot perform tasks like captioning, the closest thing to captioning they can do is give an image and a set of possible captions to select the right caption, this can be seen as a classification problem or as an image to text retrieval problem. This can be done using one of the cross-modal outputs, like the pooler output of the ViLT, the cosine similarity between text and image of CLIP, or the Cross-Modality Output of LxMERT. In each of these cases there you input the image and a possible text and obtain a score, do that for every text candidate and select the one with the highest score. To solve other key tasks like VQA and Image retrieval, a very similar process is done, also using one of these output that can be used to select the best candidate answer, let that be an image or question and an answer.

## 3.2 Vision and Language Model Decoder Architectures

Although encoding together images and text can obtain good results for tasks such as classification if we want a model that can generate text we need a decoder. Most SOTA models use Large Language Models (LLMs) as the text decoder and use the encoder only to encode the image following the problem definition we laid in chapter 1. This way in this section, we already have  $Y_v = f(I)$  using a Visual Transformer (ViT) and we want to assemble the other steps until we get our output text. With this in mind, when creating a V&L model with this kind of architecture we should ask,

- Which image tokens do we extract from the ViT?
- How does the bridging function  $h$  process the image?
- Does the bridging function  $h$  also process text?
- Does the decoder process the tokens coming from the image and the text differently?

In this section, we will start by explaining what are LLM's, see how different models answer differently to the above questions, and explain some interesting important features used to train some models.

Another important thing to consider when we are assembling these models is how we are going to train and update the weights of the encoder and the decoder, but we will focus on models that both the encoder and the decoder are frozen, i.e. their weights do not update, this is what is done by many SOTA V&L models [21, 8]. From now on we will always suppose that these pre-trained models are frozen unless explicitly stated otherwise.

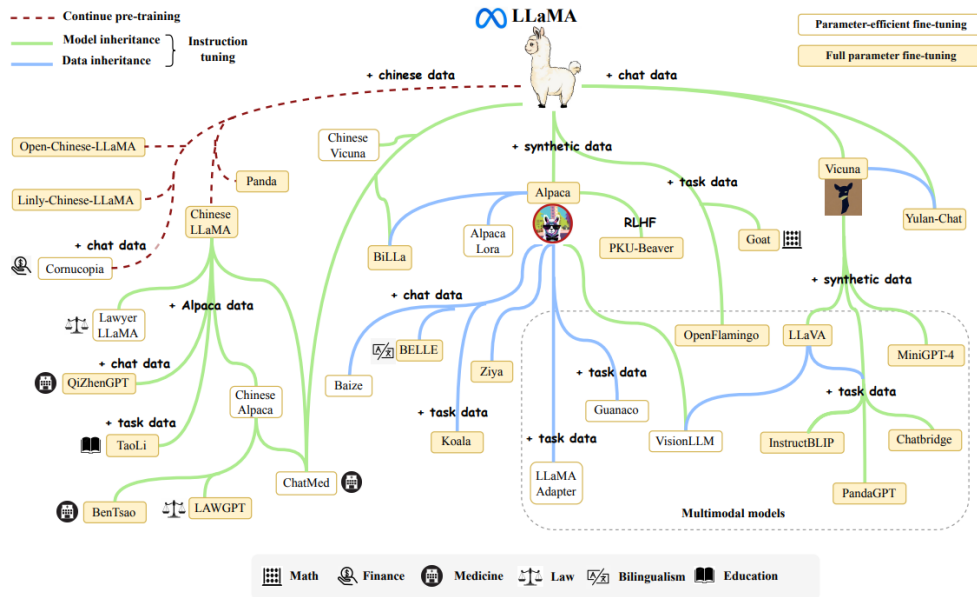


Figure 3.6: An evolutionary graph of the research work conducted on LLaMA. [46]

### 3.2.1 Large Language Models

Language Models date back to the 90s with simple models computing the probabilities of words or combinations of words from some corpus evolving to the large models we have today like GPT-4 [26]. A Language Model is a statistical model that captures word co-occurrences in a large corpus and can predict the next word in a sequence of words, according to the language data where the model was trained. Large Language Models (LLMs) are Language Models with billions of parameters trained with large amounts of data.

In order to generate open answers to questions, we need an LLM (Large Language Model). GPT-3 [3] showed that language models when trained with a lot of data can be used to solve tasks it has never seen before. With this information more Large Language Models started to appear, some commercial models like GPT-4 [26], Bard and PaLM-2 [1] but also some open models like LLaMA [36].

An advantage of these large models is that although the training is expensive and requires a lot of data, once they are trained they are general and can be fine-tuned with a lower cost for specific tasks. This is the case of Meta’s open model, LLaMa [36], which since came out a lot of research has been done based on it. Notable mentions are Alpaca [34] which was fine-tuned with 52K instruction-following demonstrations and Vicuna [6] which was fine-tuned with user-shared conversations. More examples of models built on LLaMa can be seen in figure 3.6.

### 3.2.2 ViT Token extraction

As explained in section 2.1.2.2, when an image is processed by a ViT, the former will take as input  $N$  patch tokens from the image and an extra class token. This way the last layer of this transformer will have output the embeddings  $O_v \in \mathbb{R}^{(N_p+1) \times d_v}$ , where  $N_p$  is the number of patches (not including the CLS token). But we also noted that usually the embedding of the class (CLS) token is used to represent the image, with this in mind we can get our encoder output vector  $Y_v$ , from:

- **CLS** - we can use only use the CLS token,  $Y_v \in \mathbb{R}^{1 \times d_v}$
- **All tokens** - we can use all the tokens,  $Y_v \in \mathbb{R}^{(N_p+1) \times d_v}$
- **All except CLS** - we can use all tokens except the CLS  $Y_v \in \mathbb{R}^{N_p \times d_v}$

V&L models like Fromage [14] use only the CLS, BLIP-2 [17] and Instruct-BLIP [8] use all the tokens, LLaVA [21] uses all except the CLS.

It is also important to consider from which layer to extract these embeddings, although it is common to extract from the last layer, it is also common to extract from the second to last layer because earlier layers in transformer-based models, especially the second-to-last layer, contain richer or more generalized feature representations. The argument is that this is because the last layer is often more task-specific, particularly when the model is fine-tuned for a particular downstream task. In contrast, representations from the second-to-last layer are considered more adaptable and less specialized, making them useful for a wider variety of tasks. While Fromage [14] extracts the CLS token from the last layer, most approaches that extract multiple tokens like LLaVA [21] and Instruct-BLIP [8] extract these tokens from the second to last layer.

### 3.2.3 Bridging function

After choosing what tokens to extract from ViT, we have the embeddings  $Y_v$  and we want to transform them into new embeddings  $Z_v$  that will be part of the input of the decoder. This means that now we search for our function  $h$ , here we present some of the most used functions.

#### 3.2.3.1 Linear Layer

One of the simplest and most effective functions is a linear layer. In this case, we will have a weight matrix  $W \in \mathbb{R}^{d_v \times d_t}$  that we multiply with our visual embeddings  $Y_v \in \mathbb{R}^{N \times d_v}$  and a bias component  $B \in \mathbb{R}^{N \times d_t}$  which will be summed, obtaining,

$$Z_v = h(Y_v) = Y_v W + B \in \mathbb{R}^{N \times d_t} \quad (3.7)$$

where  $B = \mathbf{1}b$ ,  $b \in \mathbb{R}^{1 \times d_t}$  and  $\mathbf{1} = [1, \dots, 1] \in \mathbb{R}^{N \times 1}$ , this means that each row of  $Z_v$  is  $z_v^r = y_v^r W + b$  where  $z_v^r$  and  $y_v^r$  are the row  $r$  of  $Z$  and  $Y$  To simplify we define a linear

layer as a function  $LL$  defined by

$$LL(Y) = YW + B \quad (3.8)$$

This approach is used by LLaVA [21] as can be seen in figure 3.7, as we mentioned earlier this model uses all the token embeddings except the CLS, meaning we get  $Z_v \in \mathbb{R}^{N_p \times d_t}$ . Using this simple linear transformation has the advantage of being very light, allowing for a large amount of data while training for a small amount of time.

Another model that uses this transformation is Fromage [14], as we can see in figure 5.1 which only uses the embedding CLS token output. This means that for the standard Fromage model, we get  $Z_v \in \mathbb{R}^{1 \times d_t}$ , being just one vector. However, this model allows for a more flexible use of the linear transformation, allowing  $h$  to transform  $Y_v$  into a sequence of  $k$  vectors. For this we have that  $W \in \mathbb{R}^{d_v \times (d_t k)}$  and  $b \in \mathbb{R}^{1 \times (d_t k)}$ , this way we have

$$Z_v^* = Y_v W + b \in \mathbb{R}^{1 \times d_v k}, \quad (3.9)$$

$$Z_v = \text{reshape}(Z_v^*) \in \mathbb{R}^{k \times d_v}, \quad (3.10)$$

therefore,

$$Z_v = h(Y_v) = \text{reshape}(Y_v W + b) \quad (3.11)$$

this means that multiplying by  $W$  gets a high-dimension vector that is then broken into a sequence of vectors. This is mathematically equivalent to having  $k$  different weight matrices each one creating one of the  $k$  vectors of the sequence.

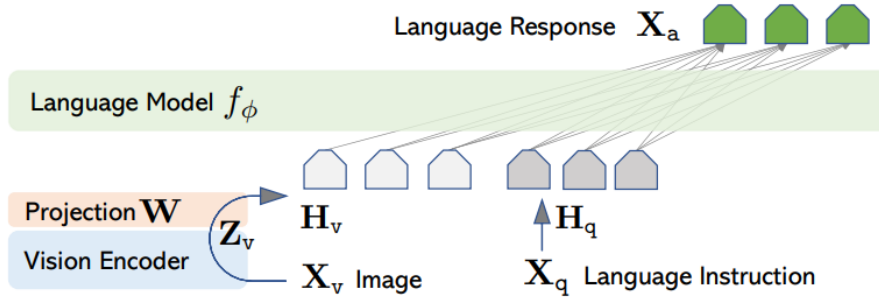


Figure 3.7: Architecture of the LLaVA model. The image here represented as  $X_v$  is transformed to  $Z_v$  by the vision Encoder entering the LLM as  $H_v$  after the linear projection. [21]

### 3.2.3.2 Multi-Layer Perceptron

A slightly more complex transformation that gets better results is to use a simple neural network with more than 1 layer, that is, a Multi-Layer Perceptron (MLP), this approach is used by LLaVA 1.5 [19] which uses 2 layers. Following this model and for the sake of simplicity we will detail how an MLP with 2 layers works. In this case, we would have 2 ( $L$  if we had  $L$  layers) weights matrices  $W_0 \in \mathbb{R}^{d_v \times d_k}$ ,  $B_0 = \mathbf{1}b_0 \in \mathbb{R}^{N \times d_k}$ ,  $W_1 \in \mathbb{R}^{d_k \times d_t}$  and

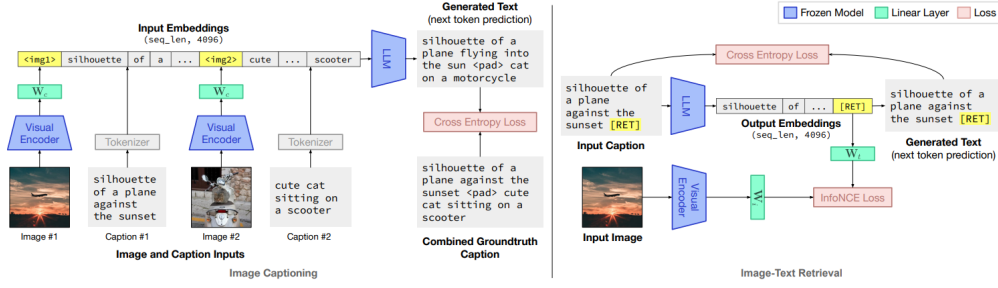


Figure 3.8: Architecture and training of the Fromage Model. The output of the Visual encoders is transformed by the linear layer  $W_c$  when entering the LLM and transformed by the linear layer  $W_i$  when used for retrieval, being then compared with the transformation of the vector [RET] by the layer  $W_t$ . [14]

$B_1 = \mathbb{1}b_1 \in \mathbb{R}^{N \times d_t}$ , obtaining:

$$Z_v^* = LL_0(Y_v) = Y_v W_0 + B_0, \quad (3.12)$$

$$Z_v^{**} = \text{activation}(Z_v^*), \quad (3.13)$$

$$Z_v = LL_1(Z_v^{**}) = Z_v^{**} W_1 + B_1 \quad (3.14)$$

that is,

$$Z_v = h(X_v) = LL_1(\text{activation}(LL_0(Y_v))), \quad (3.15)$$

where the activation is an activation function.

### 3.2.3.3 Q-Former

A more complex transformation is the one used by BLIP [17] and Instruct-BLIP [8] which consists of a Query Transformer that learns query embeddings that will interact with the image encoding through cross-attention as we can see in figure 3.9. The advantage of this strategy is that it can obtain a better transformation by using a nonlinear one, but at the cost of longer training although much lighter than without the frozen encoder and decoder.

In the former transformations, the number of embeddings in the output of  $h$ ,  $Z_v$ , was dependent on the number of embeddings in the inputs,  $Y_v$ , being the same number or multiple. Using the Q-former the number of embeddings can be reduced, being the number of output embeddings the same as the number of queries. This allows us to use all tokens but output only the amount of tokens we find useful. We can write this as:

$$Z_v = h(Y_v) = \text{Q-Former}(Y_v) \in \mathbb{R}^{N \times d_t}, \quad (3.16)$$

where  $N$  is the number of queries.

### 3.2.4 Multimodal bridging function

Although we are focusing on models where the bridging function transforms vision embeddings, there are model where the function  $h$  can also take text as input, that means

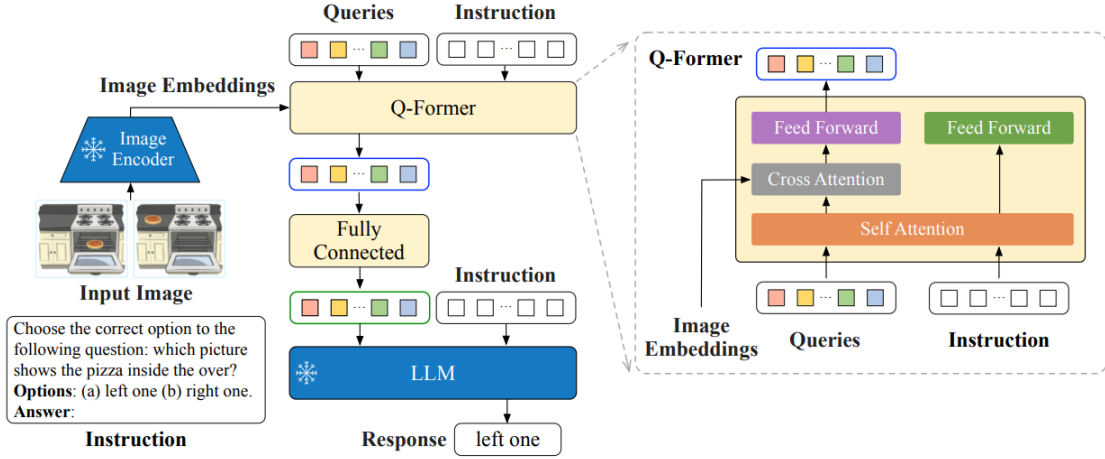


Figure 3.9: InstructBLIP [8]

that instead of

$$Z_v = h(Y_v), \quad (3.17)$$

we have

$$Z_v = h(Y_v, Y_t), \quad (3.18)$$

where  $Y_t$  is the preprocessed text we get after transforming the text tokens into vectors.

This approach is used in Instruct-BLIP as we can observe in figure 3.9 where "Instruction" tokens enter the Q-Former together with the queries. This way, we can update equation 3.2.3.3 as:

$$Z_v = h(Y_v, Y_t) = \text{Q-Former}(Y_v, Y_t) \in \mathbb{R}^{N \times d_t} \quad (3.19)$$

This approach seems to be very useful since it allows the model to extract features from the image that are relevant to the question or instruction the model is answering.

### 3.2.5 Text and Image in the Decoder

In the previously mentioned V&L models, both text and image tokens are processed by the same decoder. However, CogVLM [39] takes a different approach by treating the image and text tokens differently within the decoder to better align the frozen pre-trained models (vision encoder and language model). Specifically, CogVLM introduces a trainable visual expert module that operates on the image tokens while keeping the rest of the language model (LLM) frozen. In CogVLM, the visual expert module uses different trainable Query, Key, and Value matrices (QKV) specifically for image tokens. This allows the model to handle the image tokens with specialized attention mechanisms that differ from those used for the text tokens. The text tokens, on the other hand, continue to use the frozen QKV matrices from the pre-trained language model. In this architecture, the parts of the transformer model that process text remain entirely frozen, preserving the original language modeling capabilities of the LLM. At the same time, the introduction

of a trainable visual expert module allows the model to effectively incorporate image information by optimizing the image processing pathways.

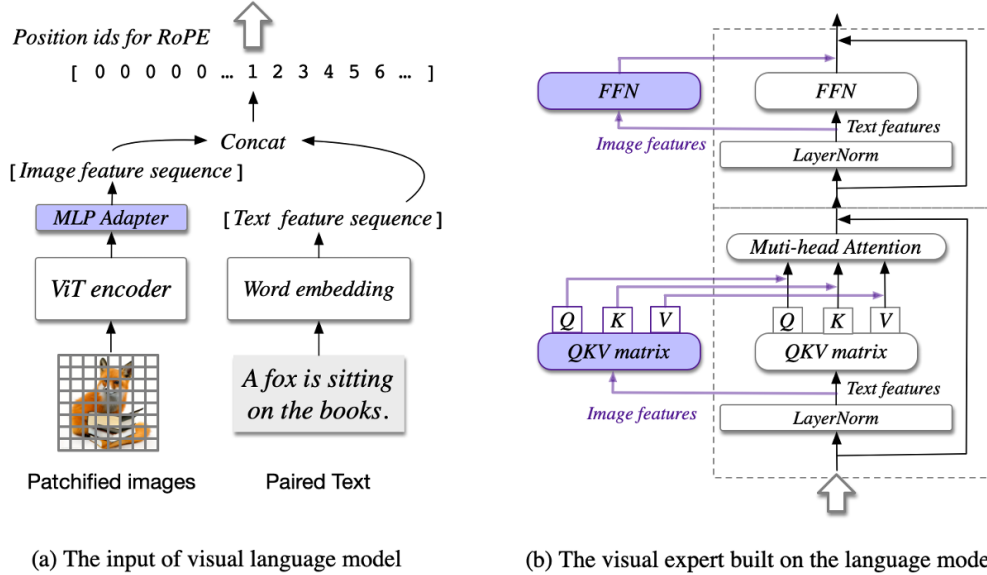


Figure 3.10: The architecture of CogVLM. (a) The illustration about the input, where an image is processed by a pre-trained ViT and mapped into the same space as the text features. (b) The Transformer block in the language model. The image features have a different QKV matrix and FFN. Only the purple parts are trainable.

### 3.2.6 Extra Vision-Language (VL) Features

Having explored the components that enable text decoders to process image inputs in vision-language models, we now turn our attention to some additional training features that enhance the performance and flexibility of these models. Two key features are instruction tuning and the use of retrieval tokens. These features allow VL models to handle a wider variety of tasks and expand their capability to include both image retrieval and complex multimodal reasoning.

#### 3.2.6.1 Instruction-Based Vision Language Models

Instruction tuning has been a significant development in natural language processing (NLP), particularly with models like InstructGPT [40] and FLAN-T5 [7]. By fine-tuning large language models (LLMs) on a variety of tasks framed as natural language instructions, the models become more adaptable and capable of following arbitrary instructions across different contexts. This approach allows models to perform diverse tasks with the same architecture, simply by presenting them with different instructions. This flexibility is now being extended to vision-language models.

In models like LLaVA, instruction tuning is done using datasets generated by text-only GPT-4. The training data includes images with captions and bounding boxes, which are

used as input to generate several tasks such as question-answer pairs, instructions, and visual reasoning scenarios. This data helps the model generalize across different tasks and prompts, making it capable of solving new vision-language tasks it has not seen before. The instructions consist of text informing the model of the task such as *"Please provide a short depiction of the picture."* and *"Briefly describe the content of the image."* [8] for image captioning tasks and *"Based on the image, respond to this question with a short answer: Question. Answer:"* [8] for VQA tasks.

### 3.2.6.2 Fromage [RET] token

The FROMAGE [14] model introduces another unique feature: the [RET] token, designed for image retrieval. Most vision-language models generate text as output, but FROMAGE extends its capabilities by allowing it to retrieve and output images as well. This capability is made possible through the use of the [RET] token, which serves as a signal for the model to retrieve images from a database.

Here's how the [RET] token works in the FROMAGE model:

- The [RET] token is a special token added to the model's vocabulary. During training, it is treated like any other token and appears within the generated text sequence. The model learns the appropriate conditions under which to output the [RET] token during training.
- When the model outputs the [RET] token during inference, this token is mapped into a retrieval embedding space. The token is processed through a projection matrix,  $W_i$ , which maps it into the same vector space as image embeddings.
- The stored image embeddings (from the model's image database) are also mapped into the same retrieval space using a matrix,  $W_i$ . Once the [RET] token embedding is generated, it is compared to the stored image embeddings, and the most similar image is retrieved based on cosine similarity.
- This enables the model to interleave text generation with image retrieval, allowing it to flexibly switch between textual and visual outputs. For example, the model might output descriptive text and then include an image that corresponds to part of the generated text, seamlessly blending retrieval into the narrative.

The [RET] token is trained together with text generation tasks by appending the token to captions during the training phase. The model attends to all tokens, including the [RET] token, and learns to use it appropriately based on the context. This additional retrieval capability is useful for tasks that involve multimodal outputs, such as suggesting relevant images for a given text or retrieving visual references during conversation.

### 3.3 Summary

In this section, we explored the architectures of some Vision and Language models focusing on two parts, Encoder Architectures and Decoder Architectures. In the Encoder Architectures, we explored how image and text can be encoded together following a unified mathematical framework that aims to encapsulate most of V&L Encoder architectures. We conclude that encoder methods are good for some tasks like classification or to use as the initial embedding of some models, but they cannot generate text leading to them not being enough for those tasks where we want generated text. In the Decoder Architectures, explored models that solve this problem. We mostly explored models that use frozen image encoders and frozen text decoders, diving into the different ways to connect them, and presenting some of the SOTA Vision and Language models.

During the writing of this thesis, other SOTA models came out, some important mentions are LLaVA-NeXT [20] following an architecture very close to LLaVA [21] but improved with some techniques such as dynamic high resolution and data mixture, and Ovis [25] which introduces to the visual encoder a visual embedding table which makes the visual patch being mapped to a probabilistic token which helps select multiple embeddings from the embedding table and output their weighted combination.

## V&L TASKS AND DATASETS

In this chapter, we explain the datasets and tasks we used to train and evaluate the models we present in the next chapters.

### 4.1 Datasets

#### 4.1.1 MSCOCO

The Microsoft Common Objects in COntext (MSCOCO [18]) dataset is a large-scale dataset designed for object detection, segmentation, and image captioning. It contains complex, everyday scenes that feature common objects in their natural contexts, making it both comprehensive and versatile. For the purpose of this work, we focus specifically on the image captioning aspect.

The dataset comprises 330,000 images, each paired with five human-generated captions. These images were sourced from Flickr by searching for combinations of 80 object categories and various scene types to capture multiple objects in realistic settings. The captions were generated by independent annotators, who adhered to specific guidelines designed to ensure consistency and clarity. These guidelines included:

- Describe all the important parts of the scene.
- Do not start the sentences with “There is.
- Do not describe unimportant details.
- Do not describe things that might have happened in the future or past.
- Do not describe what a person might say.
- Do not give people proper names.
- The sentences should contain at least 8 words.

In figure 4.1 we can observe examples of images and captions present in this dataset



Figure 4.1: Example images and captions from the Microsoft COCO Caption dataset. [18]

### 4.1.2 VisDial

Visual Dialog (VisDial), introduced by Das et al. (2017), is a large-scale dataset designed to facilitate research on visual dialog systems. It builds upon the tasks of image captioning and visual question answering (VQA) by incorporating a dialog element, where an AI agent must hold a meaningful conversation about a given image. This dataset is particularly well-suited for developing models that understand both visual content and natural language interactions.

The VisDial dataset contains over 120,000 images sourced from the MSCOCO dataset, each paired with a caption and a dialog comprising 10 rounds of questions and answers. To obtain these conversations, two workers were paired to chat in real-time. One of the workers (the 'questioner') sees only the caption of a given image. while the second worker (the 'answerer') was given both the caption and the image. The first worker's goal is to ask questions to 'image the scene better' while the second worker answers. In figure 4.2 we can see a Visual Dialog example.

### 4.1.3 CC3M

Google's Conceptual Captions (CC3M) is a dataset composed of about 3.3M image/caption pairs. Each pair contains a caption and a URL to an image on the internet, some of these images were not available anymore and others could not be opened, we ended up with about 2.7M images leading to a loss of available data of about 18.4 %, the complete numbers can be seen in the following table:

The Conceptual Captions dataset was created using a Flume [5] pipeline. This pipeline processes billions of Internet webpages in parallel. It extracts, filters, and processes candidate image-caption pairs from these webpages. This pipeline allows for the selection and transformation of the image-caption pairs, selecting the really relevant ones, this transforms the alt-text into a caption that in principle could be written only by watching



Figure 4.2: Example of a Visual Dialog [9]

Table 4.1: Comparison between the original number of images and the actually available images in the cc3m [31] dataset.

	train	validation	total
original	3318333	15840	3334173
available	2709383	12702	2722085
loss	18.4%	19.8%	18.4%



Figure 4.3: Example of images, alt-text and conceptual caption from the CC3M dataset [31].

at the image, in figure 4.3 we can compare the original alt-text with the new caption used by CC3M.

Even with this transformation, this dataset's captions are not always good since it is automated and not human-made.

## 4.2 Tasks

### 4.2.1 Captioning

The image captioning task involves generating a descriptive text (caption) for a given image. It combines both computer vision and natural language processing to understand the visual content of the image and describe it in human language. We can formally define

it as:

$$\Phi_{caption}(I) = c \quad (4.1)$$

where  $I$  is an image and  $c$  is the generated caption, we can also add a prompting instruction  $P$  that would condition the answer on things like the focus or the level of detail., obtaining:

$$\Phi_{caption}(I, P) = c \quad (4.2)$$

When we want to evaluate the image captioning performance of a certain model using a dataset such as MSCOCO we can compare the ground-truth answers from the dataset with the ones generated by the model. This comparison is made using some metrics like the following:

**BLEU-n** (Bilingual Evaluation Understudy) evaluates the precision of n-grams (sequences of n words) between the generated and reference texts. The "n" refers to the size of the n-grams being compared:

- BLEU-1 measures unigram precision (individual words),
- BLEU-2 measures bigram precision (pairs of words), and so on.
- BLEU-n combines the precision of multiple n-gram sizes, typically up to n=4 (BLEU-4), with a brevity penalty to avoid rewarding overly short outputs. Scores range from 0 to 1, with higher values indicating a closer match. BLEU focuses on precision, not recall, and doesn't account for synonymy or meaning.

**METEOR** (Metric for Evaluation of Translation with Explicit ORdering) is designed to improve upon BLEU by incorporating both precision and recall, while also accounting for synonymy and stemming. It aligns words between the generated and reference texts and rewards matches based on exact words, synonyms, and root forms. METEOR values range from 0 to 1, with higher scores indicating better matches, and it also includes a penalty for fragmented matches.

**CIDEr** (Consensus-based Image Description Evaluation) is specifically tailored for evaluating image captioning tasks. It measures the consensus between the generated caption and a set of reference captions, with an emphasis on content words. CIDEr uses TF-IDF (term frequency-inverse document frequency) to give higher weight to important and unique words, focusing on how well the generated caption captures key visual elements. Scores range from 0 to 1, where higher values indicate a better match with the reference captions.

## 4.2.2 Retrieval

Image retrieval is a task whose main goal is to find relevant images from a dataset based on a given input query, typically in the form of text. This task combines computer vision

and natural language processing to understand the visual content of each image and determine its relevance to the query. We can formally define it as:

$$\Phi_{retrieval}(q) = \mathbf{I} \quad (4.3)$$

where  $q$  is the input query, and  $\mathbf{I}$  represents the set of retrieved images. The most common approach for image retrieval is to rank images based on their similarity to the query. This allows the system to select the top  $k$  most relevant images from the dataset, depending on how many images we want to retrieve.

Image retrieval can also be evaluated in tasks using datasets like MSCOCO, where the query is a caption, and the goal is to retrieve the corresponding image. Several metrics are used to evaluate retrieval performance, including:

- **Precision** =  $\frac{\text{Relevant retrieved images}}{\text{All retrieved images}}$ : This measures the proportion of retrieved images that are relevant.
- **Recall** =  $\frac{\text{Relevant retrieved images}}{\text{All relevant images}}$ : This evaluates the proportion of all relevant images that are successfully retrieved.

Since it is common to use a ranking, we can also use metrics like precision@k and recall@k where the retrieved images are the top k ranked images for the given query.

### 4.2.3 Visual Question Answering

Typically the VQA model is formulated in one of two ways: as a classification problem or as a generative problem. This gives us two different approaches the closed-answer approach and the open-answer approach. Let  $I$  be an image,  $T$  a text corresponding to a question in natural language and  $A = (a_1, a_2, \dots, a_n)$  a closed-set of  $n$  possible answers. We can represent classification Visual Question Answering models as functions  $\Phi$  that can be expressed as

$$\Phi_{closed}(I, Q, A) = a \in A. \quad (4.4)$$

For generative Visual Question Answering models, we don't have a pre-determined set of possible answers, this way they can be expressed simply as

$$\Phi_{open}(I, Q) = a, \quad (4.5)$$

where  $a$  corresponds to the generated answer, i.e. a sequence of text tokens.

Both approaches can be extended to use context, which means that instead of having just an image and a question we might also have a text that gives some context which may alter the answer. This leads to the following formulations:

$$\Phi_{closed}(I, C, Q, A) = a \in A. \quad (4.6)$$

$$\Phi_{open}(I, C, Q) = a. \quad (4.7)$$

where  $C$  is the context and  $a$  the answer given by the model.

# A PORTUGUESE VISION AND LANGUAGE MODEL

In this chapter we present a Portuguese V&L model based on FROMAGE, changing the LLMs and the dataset it was trained on.

## 5.1 Introduction

Vision and Language are two of the main communication and information perception mediums, serving as fundamental channels through which humans interpret and interact with the world around them. Devising Vision and Language (V&L) models that can seamlessly combine these two modalities is paramount to delivering AI systems capable of addressing tasks such as image captioning and visual question-answering, essential tasks to aid visually impaired individuals, and Image-to-Text and Text-to-Image retrieval, for general multimodal information seeking. Recently, there have been notable advances in vision and language models [21, 14, 12], which leverage Large Language Models as backbones [35, 3, 45] (LLMs). Most of these advances have been made with models in English or other high-resource languages, leaving behind other lower-resource languages, as is the case of European Portuguese (PT-PT). This evidences the urgent need of developing large V&L LMs, openly available, for PT-PT speakers. This raises two complementary challenges: 1) how to overcome the limited availability of PT-PT multimodal datasets and resources, and 2) how to train a Large Vision and Language model, capable of addressing multiple V&L tasks, in PT-PT.

Most LLMs are trained with text-only web scraped data, achieving great performance on a myriad of natural language tasks, but lack an overall understanding of images, thus not having visual reasoning capabilities. Pioneering vision and language models, adopted fully multimodal Transformer-based models [24, 42, 38], with either single-stream or dual-stream architectures [4], pre-trained on image-text pairs. More recently, towards generalizing high-performing large LMs to the visual domain, it is common practice to leverage text-only LLMs as the backbone and equip them with a visual encoder [28, 10].

Then, LLMs are augmented with a visual projection component that aligns visual tokens with the LLM token-space [14, 21].

In this chapter, we seek to deliver a European Portuguese vision and language LM. To this extent, we make two major contributions: **1)** we create and make available<sup>1</sup> both large-scale image-text pre-training datasets as well as well-known V&L benchmarks in European Portuguese. In particular, CC3M [31] PT-PT (3 million image-caption pairs) for pre-training, MSCOCO [18] PT-PT (image-caption pairs) and VisDial [9] PT-PT (visual dialogs) for benchmarking on downstream tasks V&L tasks. An extensive assessment of available machine translation approaches is carried out. **2)** following the V&L LMs state-of-the-art, we adapt the FROMAGe [14] model to support PT-PT. Its flexible decoder-based architecture, augmented with multimodal specialized layers, gives the model the capacity to process and produce interleaved multimodal inputs and outputs. Given that a key step is to replace the original LLM by a PT-PT LLM, we leverage a recent PT-PT text-only decoder, Glória [22], and conduct extensive experiments, in a zero-shot setting, on image caption and visual dialog tasks.

Most of existing models are in English or other high-resource languages. Very recently, open European Portuguese LLMs have been made available. In particular, Glória [22] is a European Portuguese LLM Decoder based on GPTNeo [2] - which approximates the GPT3 architecture - trained on a 35B token corpus, from a diverse set of domains, including web content, news pieces, encyclopedic knowledge, news articles, and dialogs. Gervásio [30] is another relevant European Portuguese LLM decoder which is based on a pre-trained LLaMA 2 7B [35] model, fine-tuned on Portuguese instruction datasets, comprising around 83M tokens. Regarding V&L approaches, literature is scarce. CAPIVARA [29] trains a Brazilian Portuguese CLIP model, while performing data augmentation through image captioning and machine translation. In this work, and using recent developments in the LLM PT-PT, we seek to narrow this gap, by introducing a European Portuguese V&L model.

## 5.2 EU-Portuguese Datasets for Vision and Language AI

Due to the lack of European Portuguese V&L datasets, we embraced the task of translating core vision and language datasets from English into European Portuguese. Given the size of existing datasets (millions scale), translating the datasets with human experts would be too costly, hence we considered three distinct automatic machine translation models: first, we considered a) **MADLAD-400** [15], a model trained on a 3T token dataset based on CommonCrawl, created by Google, covering text data from over 400 languages; b) **Narrativa**<sup>2</sup>, which is an mBART-50 [33] model fine-tuned on the opus-100 [44] dataset for English to Portuguese Translation, c) **DeepL**<sup>3</sup> a commercial translation service.

---

<sup>1</sup>Datasets will be made available after publication.

<sup>2</sup><https://huggingface.co/Narrativa/mbart-large-50-finetuned-opus-en-pt-translation>

<sup>3</sup><https://www.deepl.com/translator>

We started by assessing the performance of each of the three approaches. First, although MADLAD-400 seems to give good translations, most are in Brazilian Portuguese. Narrativa translations are in European Portuguese, but for many captions, the model output is the original English caption, rather than its translation. DeepL seems to solve these problems, by providing high-quality European Portuguese translations, with the drawback of being a commercial solution. Table 5.1 illustrates some of the translated examples of the CC3M dataset [31]. For example, for the first image, Narrativa outputted the original caption, in the second image MADLAD-400 uses Brazilian Portuguese lexicon in its translations (e.g. "*ônibus*" instead of "*autocarro*", the word bus). Something we also notice is that MADLAD-400 often does not translate the full caption (as in the first image).

Given these observations, we translated the CC3M [31], MSCOCO [18], and the VisDial [9] datasets, using DeepL, creating the CC3M PT-PT, MSCOCO PT-PT and VisDial PT-PT, respectively. The CC3M PT-PT dataset was used as the pre-training dataset, and both MSCOCO PT-PT and VisDial PT-PT were used for benchmarking retrieval, image-captioning and visual dialog tasks. Table 5.2 shows the aggregated statistics of these datasets. It is important to note that the lower number of total tokens in the Narrativa translation stems from the fact that some captions are not actually translated. DeepL translations have higher token numbers than the original English dataset, which despite corroborating with the increased verbosity of Portuguese vs. English, will have an impact on the models' performance.



**Original\*:** plenty of space : at square feet the property would have ample room for actor and her daughter

MADLAD-400: abundância de espaço: em pés quadrados a propriedade teria amplo espaço

Narrativa: plenty of space : at square feet the property would have ample room for actor and her daughter

DeepL: muito espaço: em metros quadrados, a propriedade teria muito espaço para o ator e a sua filha

**Original\*:** people waiting for the bus in snow storm



MADLAD-400: pessoas à espera do ônibus na tempestade de neve

Narrativa: Pessoas à espera do autocarro em tempestade de neve

DeepL: pessoas à espera do autocarro numa tempestade de neve

**Original\*:** person serves lunch to two of her daughters at their farm.



MADLAD-400: uma mulher serve o almoço para duas de suas filhas em sua fazenda

Narrativa: A pessoa serve o almoço a duas filhas da fazenda dela.

DeepL: uma pessoa serve o almoço a duas das suas filhas na sua quinta.

Table 5.1: Translation results of sample captions from the CC3M dataset, using each of the three considered translation approaches.

Table 5.2: Translation statistics, for CC3M and COCO, with different machine translation approaches. # Samples - total number of samples, # Tokens - total number of tokens, # Avg. Tokens/Sample - average number of tokens per sample. \* Stands for the original captions.

	Statistic	English*	MADLAD	Narrativa	DeepL
CC3M	# Samples	2 709 383	2 709 383	2 287 769	2 709 383
	# Tokens	27 919 393	26 558 075	24 257 997	29 844 147
	# Tokens/Sample	10.30	9.80	10.60	11.02
COCO	# Samples	25 014	25 014	23 614	25 014
	# Tokens	282 297	282 172	267 893	292 626
	# Tokens/Sample	11.29	11.28	11.34	11.70

## 5.3 Method

In this section we present V-Glória, an European Portuguese Large V&L model, capable of flexibly interleaving the two modalities, images and text, and therefore generalize to different NLP and CV tasks such as multimodal retrieval, image captioning, and visual dialog. Therefore, we adapt the FROMAGe model [14] architecture, which leverages a text LLM and adds a set of projection layers to align images with the LLM input subspace, and support generative retrieval. Specifically, it allows us to use an European Portuguese LLM, that will be frozen during training, with lightweight training strategies aimed at equip V-Glória with visual and linguistic reasoning capabilities.

### 5.3.1 V-Glória Architecture

#### 5.3.1.1 PT-PT Language Model Backbone.

V-Glória uses a Portuguese large language model decoder originally trained with text-only data with a causal language modeling task. V-Glória is based on a PT-PT open and top performing LLM, Glória [22]. In the experiments, we compare it with alternative LLM backbones, such as Gervásio [30].

#### 5.3.1.2 Visual Encoder Model

Images are encoded using a pre-trained CLIP ViT-L/14 [28], such that given an image  $y$ , the visual model outputs  $v(y) \in \mathbb{R}^m$ , corresponding to the [CLS] token embedding. Both  $\theta$  and  $\phi$ , both LLM and visual encoder parameters will be frozen.

#### 5.3.1.3 Visual Projection Layer

With the LLM and the visual encoder frozen, a projection layer is used to map the encoded images to the embedding subspace of the LLM token. Namely, a linear layer,  $v(y)^T \cdot \mathbf{W}_c \in \mathbb{R}^d$ , where  $d$  corresponds to the LLM hidden dimension.

#### 5.3.1.4 Multimodal Retrieval.

In order to support retrieving images, conditioned either on text or images, a special token [RET] is added to the model vocabulary, so that at any point in the decoding, the model can decode this token, and, its embedding (which will be learned) can be used for retrieval. During training, a [RET] token is appended to the end of the input captions. In practice, two linear mappings are trained,  $\mathbf{W}_t \in \mathbb{R}^{d \times q}$  and  $\mathbf{W}_i \in \mathbb{R}^{m \times q}$ , which will map the hidden representation of [RET] obtained from the last hidden layer of the LLM and the visual embeddings, respectively, into a common  $q$  dimensional space.

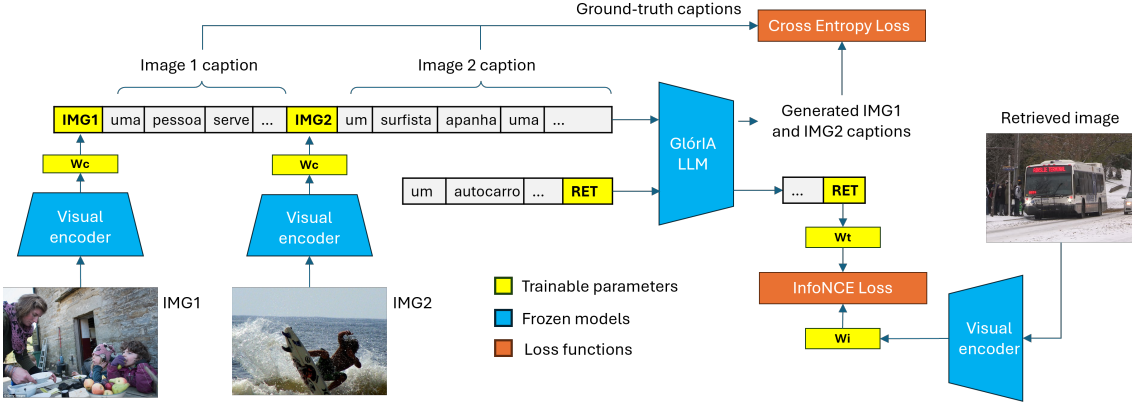


Figure 5.1: Overview of the V-Glória architecture. The model is trained on image-text pairs for image captioning and image-text retrieval. The LLM and visual encoder are frozen, while the three projection layers (in yellow), with weight matrices  $\mathbf{W}_c$ ,  $\mathbf{W}_i$ , and  $\mathbf{W}_t$ , are learned.

### 5.3.2 Training

The training tasks are specifically designed to equip the model vision and language reasoning capabilities: describing visual content; processing interleaved images and text in its context; and third matching images to text and vice-versa. The model is trained with a multi-task objective  $\mathcal{L}$  comprising image captioning and image-text retrieval, as illustrated in Figure 5.1, with

$$\mathcal{L} = \lambda_c \mathcal{L}_c + \lambda_r (\mathcal{L}_{i2t} + \mathcal{L}_{t2i}), \quad (5.1)$$

with  $\lambda_c = \lambda_r = 0.5$ .

#### 5.3.2.1 Image Captioning.

For captioning, the model is trained to autoregressively predict the next token, with a Cross-entropy loss conditioned on the image representation, using

$$l_c(x, y) = \sum_{t=1}^T \log p_{\theta}(s_t | v(y)^T \mathbf{W}_c, s_1, \dots, s_{t-1}) \quad (5.2)$$

where  $s_t$  represents the  $t$ -th token of the caption  $x$ ,  $v(y)^T$  the output of the visual encoder,  $\mathbf{W}_c$  the weights of the visual projection layer, and  $\theta$  the frozen parameters of the LLM.

#### 5.3.2.2 Image-text Retrieval

For bidirectional multimodal retrieval, given a caption  $x_i$  and its corresponding image  $y_i^4$ , the InfoNCE [infonce] loss for multimodal contrastive learning is used as

$$\mathcal{L}_{t2i} = -\frac{1}{N} \sum_{i=1}^N \left( \log \frac{\exp(x_i \cdot y_i / \tau)}{\sum_{j=1}^N \exp(x_j \cdot y_j / \tau)} \right), \quad (5.3)$$

where  $x_i \cdot y_i$  corresponds to the cosine similarity between embeddings. The loss in the opposite direction,  $\mathcal{L}_{i2t}$ , is defined reciprocally, with  $x_i$  and  $y_i$  swapped.

### 5.3.3 Implementation details

To encourage the model to attend more explicitly to images, distinct examples are concatenated together with a probability of 0.5, [14] found that this is helpful in training the model to attend to the correct image within a sequence. This model is implemented in Pytorch [27], as most of the model parameters are frozen this method is memory and compute efficient, backpropagating through the frozen LLM and a visual model, but only compute gradient updates for the 3 trainable linear layers and the [RET] embeddings. We use the Adam [13] optimizer with a learning rate of 0.0003 and warmup of 100 steps to train for 2000 steps, corresponding to 1 epoch with a batch size of 180. The visual embeddings coming from our visual encoder model have dimension  $m = 1024$ , we use  $q = 256$  as the retrieval embedding dimension and the embedding dimension  $d = 2048$  inherited from Glória.

## 5.4 Experimental Setup

We assess the performance of our model in both image retrieval and image-and-text generation tasks. The models were trained on the CC3M PT-PT dataset, originally comprising 3.3 million image-text pairs, which after filtering out missing and corrupted images resulted in a total of 2.7M samples. We consider both Glória 1.3B <sup>5</sup> and Gervásio 7B <sup>6</sup> as the PT-PT LLM backbones.

Multimodal retrieval and image captioning are evaluated in both the CC3M PT-PT (full-shot) and MSCOCO PT-PT (zero-shot) evaluation sets. Models are also evaluated in the Visual Dialog task [9], in a zero-shot setting. To establish a comparison between English and EU-Portuguese, we consider the architectural twin of Glória 1.3B, GPTNeo 1.3B [2], an English-only LLM.

<sup>4</sup>For the sake of notation simplification,  $x_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}^d$  correspond to the [RET] token output of the retrieval mapping  $\mathbf{W}_t \in \mathbb{R}^{d \times q}$ , and to the outputs of the visual mapping  $\mathbf{W}_i \in \mathbb{R}^{m \times q}$ , respectively.

<sup>5</sup><https://huggingface.co/NOVA-vision-language/GlorIA-1.3B>

<sup>6</sup><https://huggingface.co/PORTULAN/gervasio-7b-portuguese-ptpt-decoder>

Table 5.3: Retrieval results for CC3M PT-PT and MSCOCO PT-PT datasets.

	LLM	Data Language	IT2T		T2I	
			R@5	R@10	R@5	R@10
CC3M	GPT-Neo 1.3B	English	13.7	31.3	11.9	29.0
	Glória 1.3B	PT-PT MADLAD-400	22.5	44.9	22.0	44.1
	Glória 1.3B	PT-PT DeepL	<b>23.4</b>	<b>45.9</b>	<b>23.3</b>	<b>45.9</b>
	Gervásio 7B	PT-PT MADLAD-400	15.5	33.8	15.3	34.4
	Gervásio 7B	PT-PT DeepL	16.6	34.8	16.1	35.5
MSCOCO	GPT-Neo 1.3B	English	21.0	30.7	21.1	29.6
	Glória 1.3B	PT-PT MADLAD-400	<b>34.7</b>	<b>46.8</b>	<b>35.7</b>	<b>47.2</b>
	Glória 1.3B	PT-PT DeepL	30.2	41.1	30.1	40.7
	Gervásio 7B	PT-PT MADLAD-400	16.6	25.5	16.5	24.2
	Gervásio 7B	PT-PT DeepL	16.7	24.5	15.7	22.3

## 5.5 Results and Discussion

In this section, we discuss the experimental results in the image captioning and cross-modal retrieval tasks. We start by evaluating our model in cross-modal retrieval, in both Image to Text (I2T) and Text to Image (T2I) settings, and then in image captioning. We follow related work, and for cross-modal retrieval experiments adopt as metrics Recall@5 (R@5), and Recall@10 (R@10), and for image captioning BLEU and METEOR. Finally, we consider the challenging task of Visual Dialog, in a zero-shot setting. For the three tasks, we follow the established task protocols.

### 5.5.1 Cross-Modal Retrieval Results

We can observe in Table 5.3 that V-Glória trained, using Glória as LLM, and with data translated with DeepL, has the best results, significantly outperforming Gervásio in both directions, although the latter has more than five times the number of Glória parameters. In the MSCOCO validation set (unseen data), we observe a similar trend, where Glória shows to be preferable to Gervásio. However, in MSCOCO, we observe that higher performance is achieved when training and evaluating using the dataset translations obtained with MADLAD-400. When comparing the performance between the two languages, i.e. PT-PT (Glória 1.3B) and English (GPT-Neo 1.3B), we observe that performance is higher in PT-PT. This shows the robustness of our training procedure and hints at the promising capabilities of PT-PT vision and language models.

### 5.5.2 Image Captioning Results

Table 5.4 shows the image captioning results. First, we observe the same trend in which our model, V-Glória, using Glória 1.3B as its LLM backbone, consistently achieves superior

Table 5.4: Captioning results on the validation split of the CC3M PT-PT and MSCOCO PT-PT datasets.

	LLM	Data Language	BLEU1	BLEU2	BLEU3	BLEU4	METEOR
CC3M	GPT-Neo 1.3B	English	18.5	9.9	6.0	4.0	17.6
	Glória 1.3B	PT-PT MADLAD-400	11.9	6.0	3.5	2.3	13.9
	Glória 1.3B	PT-PT DeepL	11.8	5.7	3.3	2.2	13.7
	Gervásio 7B	PT-PT MADLAD-400	9.6	5.4	3.4	2.3	12.3
	Gervásio 7B	PT-PT DeepL	10.8	6.1	3.8	2.6	13.1
MSCOCO	GPT-Neo 1.3B	English	42.8	24.1	12.9	7.0	13.1
	Glória 1.3B	PT-PT MADLAD-400	29.7	16.2	8.8	4.7	13.8
	Glória 1.3B	PT-PT DeepL	25.8	12.7	6.8	3.6	12.1
	Gervásio 7B	PT-PT MADLAD-400	21.6	12.3	7.0	3.9	13.4
	Gervásio 7B	PT-PT DeepL	23.8	13.3	7.9	4.7	12.9

Table 5.5: Zero-shot results on VisDial [9], for image-and-text-to-text (IT2T) and text-to-image (T2I) retrieval. Unlike previous methods, is capable of generating free-form text interleaved with image outputs through text-to-image retrieval.

LLM Backbone	IT2T		T2I	
	R@5	R@10	R@5	R@10
Glória 1.3B	4.2	14.1	17.3	25.2
Gervásio 7B	4.0	13.9	8.2	14.0

performance, compared to the Gervásio LLM backbone. Second, it can be seen that the task is much more challenging on CC3M-PT, with all models obtaining a lower performance. These low BLEU scores on CC3M, might be explained by the fact that since CC3M captions are collected from the web, and not manually annotated like in MSCOCO, making them more prone being unaligned with the image. ( see first example of Table 5.1, where the caption mentions "actor and her daughter" which cannot be guessed from the picture). However, in MSCOCO, higher BLEU and METEOR scores are obtained. It should be noted that for MSCOCO, models are evaluated in a zero-shot setting, evidencing that V-Glória is capable of generalizing to unseen data.

When comparing a full English setup vs. a PT-PT model trained on PT-PT data, we observe that the former achieves higher performance in both datasets. Given the proximity of the image captioning task to the original LLM loss, and the fact that GPT-Neo 1.3B was pre-trained on a significantly larger text corpus, compared to Glória and Gervásio, this is not surprising and we believe that this can be countered with an improved PT-PT LLM.

### 5.5.3 Visual Dialog Results

To assess our model performance on a more challenging vision and language task, we evaluate it on the Visual Dialog (VisDial) [9] task, in zero-shot, in two different settings: **a)** IT2T (image and text to text) where given an image, a dialog about it, and a question, the

model has to select the correct answer from a pool of 100 candidate answers, and **b)** T2I (text to image), where given a dialog about an image, the model has to retrieve the correct image. Given that V-Glória is an autoregressive decoder, we follow the protocol of [14] for IT2T, and given a question and answer sequence, we select the answer with the lowest perplexity, among the candidate answer options.

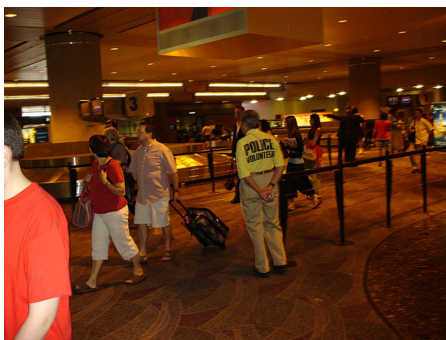
Table 5.5 shows the results. We observe that all models exhibit low performance, regardless of the PT-PT LLM backbone. Performance is, however, higher in T2I, compared to IT2T, which is consistent with the fact that the T2I task is closer to the vision and language tasks considered in training. Notwithstanding, V-Glória, using the Glória PT-PT LLM, demonstrates better generalization capabilities to new tasks, significantly outperforming the model using the Gervásio PT-PT LLM. We believe that part of these results can be dramatically improved by using a stronger PT-PT LLM. That is, despite the higher effectiveness of the Glória PT-PT LLM, it was not trained on instructions. This makes the model struggle when instructed to answer questions.

**Query:** "Uma moto Honda preta estacionada em frente a uma garagem."

**Retrieved images:**



(a) Image retrieval



**Ground truth caption:**

Várias pessoas caminham pelo aeroporto enquanto esperam pelas suas malas.

**V-Glória generated caption:**

A fila de pessoas que se encontram a caminho do aeroporto.

**Gervásio generated caption:**

pessoas a descer a passagem de nível.

(b) Image captioning.

**Question:** Quantas pessoas estão na foto?

**Answer (GT):** 5

**Answer (V-Glória):** 13

**Question:** Estão virados para a câmara?

**Answer (GT):** sim

**Answer (V-Glória):** sim

**Question:** Estão a usar casacos?

**Answer (GT):** sim

**Answer (V-Glória):** sim

**Question:** Existem árvores visíveis?

**Answer (GT):** sim

**Answer (V-Glória):** branco



(c) Visual dialog.

Figure 5.2: V-Glória can solve core vision and language tasks.

## 5.6 Conclusions

In this chapter, we introduced V-Glória, the first European-Portuguese Vision and Language model, capable of addressing multimodal tasks such as retrieval, image captioning, and visual dialogs illustrated in Figure 5.2. In addition, we will release the PT-PT high-quality translations of the most popular V&L datasets to foster research in this area. Experiments, leveraging current best performing open PT-PT LLMs as backbones, reveal performances that are competitive with the English counterparts on these tasks. V-Glória demonstrated to be capable of generalizing to unseen data, especially in multimodal retrieval. For more challenging tasks, such as Visual Dialog, the proposed approach is still not on par with English models. However, we believe that as better PT-PT models arise, including instruction-tuned ones, the performance gap can be narrowed down with the adopted methodology and leveraging our contributed PT-PT data resources.

## ABLATION STUDY: CLS BASED MODELS

In this chapter, we present an ablation study for Vision and Language models following the FROMAGE architecture using CLS tokens as image embeddings. For this analysis, we use CLIP ViT-L/14 [28] as the image encoder which is the image encoder used by FROMAGE and we chose Gemma-2B instruct version as it is a very light model while still performing well, being of the best LLMs with only 2B parameters, allowing us to train many models as required by an ablation study while not taking too long to train the models. We follow the original FROMAGE training and details, training with a batch size of 180 for 20000 steps (while they state using 18000, we opted for 20000 as it was the default setting), while their original plan was to train for 1 epoch, our training with 20000 steps and less training examples because of the loss reported in section 4.1.3 becomes slightly more than 1 epoch.

We evaluate these models with 3 different datasets. For the CC3M and MSCOCO datasets, we present the validation metrics for captioning and retrieval, it is important to note that although the models were trained using CC3M, the results for MSCOCO are more interesting since it is a dataset the model has not seen during training but also because its captions are human-made and can be inferred from the image while most image-caption pairs of CC3M are of poor quality for evaluation. We also present examples of captions generated by the models from the MSCOCO dataset images. For the VQA dataset, we present the answer accuracy results for the type of answer (yes/no, number, other). In this dataset we used the approach used by FROMAGE to evaluate it where it prunes the number of output words to the number of the longest answer and following the official VQA evaluation code. This may lead to some problems underestimating the VQA capabilities of some models by classifying as wrong an answer that is right. We also present examples of VQA answers of these models that may help to better compare models. Because of that we focus on these examples and on the other results, mainly of MSCOCO captioning.

It is also important to note that while these models are capable of both producing text from image or from image and text and also of retrieving images from text or from images and text, VQA only evaluates the former. Since training these models is a stochastic

process, where we obtain different results each time they are trained, we train each model 3 times and the results we present are the mean of these runs results.

## 6.1 ViT Layer

We have settled that the function  $f$  corresponds to the output of the image encoder and that our encoder is a CLIP ViT-L/14 [28] but that does not mean that our function  $f$  is not already defined because as explained earlier we can use the default output, that is the embedding of the CLS token in the last layer but we can choose different layers. With this in mind, we define  $ViT_{CLS}^l$  as the function that extracts the CLS token embeddings of the layer  $l$ .

The original Fromage model extracts the visual embeddings from the last layer,  $L$ . This means that in the original model

$$f(I) = ViT_{CLS}^L(I), \quad (6.1)$$

but as explained previously there are models like LLaVA [21] which extract from the second to last layer. Although they state that using this second to last layer is better LLaVA also uses all tokens instead of just the CLS therefore their conclusion may not generalize to CLS-based models. With this in mind, we start by evaluating this model with the CLS token embedding extracted from different layers.

We experiment for the 4 last layers, meaning that our function  $f$  will be:

$$Y_v = f(I) = ViT_{CLS}^l, l \in \{L-3, L-2, L-1, L\} \quad (6.2)$$

where  $L$  is the last layer and therefore  $L-3$  is the fourth last layer, with the model we used this means that  $Y_v \in \mathbb{R}^{1 \times 1024}$ . In these four models, we are using a linear layer as the bridging function  $h$ , recalling section 3.2.3.1 this means that

$$Z_v = h(Y_v) = Y_v W + b \in \mathbb{R}^{1 \times 2048} \quad (6.3)$$

where  $W \in \mathbb{R}^{1024 \times 2048}$  and  $b \in \mathbb{R}^{1 \times 2048}$  are the trainable parameters. This approach follows is illustrated in figure 6.1

Table 6.1: CC3M validation metrics for different layers.

Layer	Captioning					I2T Retrieval		T2I Retrieval	
	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	r1	r5	r1	r5
L	16.30	8.23	4.57	<b>2.74</b>	13.83	<b>23.30</b>	<b>46.02</b>	<b>22.32</b>	<b>44.41</b>
L-1	<b>16.78</b>	<b>8.45</b>	<b>4.64</b>	<b>2.74</b>	<b>13.86</b>	21.50	43.76	20.53	42.67
L-2	13.98	7.04	3.84	2.22	11.66	19.67	40.92	18.36	39.38
L-3	14.41	7.28	4.01	2.35	11.66	16.66	37.58	15.07	34.51

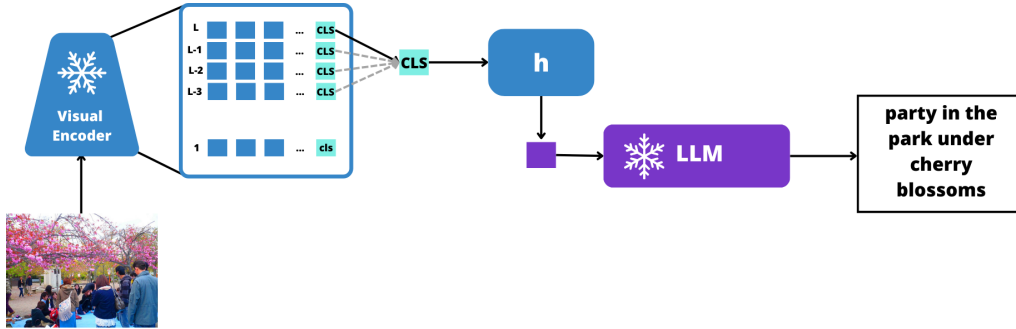


Figure 6.1: Architecture of our model generative part where  $h$  is a linear layer that transforms a CLS token, where we change the layer we extract the token. The solid black line corresponds to the layer  $L$  used by the original model, but we try for the last 4 layers.

Table 6.2: VQA validation results for different layers.

Layer	overall	yes/no	number	other
L	<b>31.28</b>	62.36	<b>3.00</b>	<b>15.13</b>
L-1	30.53	61.95	1.94	14.20
L-2	29.72	<b>62.43</b>	2.14	12.13
L-3	28.64	59.07	1.84	12.58




Table 6.3: MSCOCO captioning validation metrics for different layers.

Layer	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr
L	18.87	9.33	4.48	2.22	7.68	21.09	20.64
L-1	23.66	12.38	6.21	3.20	<b>9.20</b>	<b>23.12</b>	<b>22.50</b>
L-2	18.78	10.69	5.86	3.27	6.71	17.08	17.72
L-3	<b>30.45</b>	<b>16.81</b>	<b>8.89</b>	<b>4.74</b>	9.09	21.72	19.80

Table 6.4: MSCOCO retrieval validation metrics for different layers.

Layer	I2T		T2I	
	r1	r5	r1	r5
L	<b>14.99</b>	<b>33.97</b>	<b>15.26</b>	<b>34.89</b>
L-1	14.29	33.54	14.22	33.58
L-2	11.91	30.59	13.00	31.17
L-3	11.25	28.56	10.63	27.61

Table 6.5: VQA answers generated by the models extracting the cls embedding from different layers.

		
<p><b>Question:</b> What is the man doing in the street?</p>	<p><b>Question:</b> What are the people in the background doing?</p>	<p><b>Question:</b> What color is the building?</p>
<p><b>Layer L:</b> he is walking down the street.</p>	<p><b>Layer L:</b> they are watching a skateboarder perform.</p>	<p><b>Layer L:</b> gray.</p>
<p><b>Layer L - 1:</b> he is walking to work.</p>	<p><b>Layer L - 1:</b> the</p>	<p><b>Layer L - 1:</b> gray.</p>
<p><b>Layer L - 2:</b> he's selling ice cream.</p>	<p><b>Layer L - 2:</b> they are watching the skateboarder perform.</p>	<p><b>Layer L - 2:</b> the building is painted in a light beige color.</p>
<p><b>Layer L - 3:</b> he is driving a car .</p>	<p><b>Layer L - 3:</b> they are watching the skateboarder perform.</p>	<p><b>Layer L - 3:</b> the building is a light beige color .</p>



*L*: The cat is sitting in the car.

*L* – 1: a cat in a car.

*L* – 2: a cat in a car.

*L* – 3: a cat in a car.

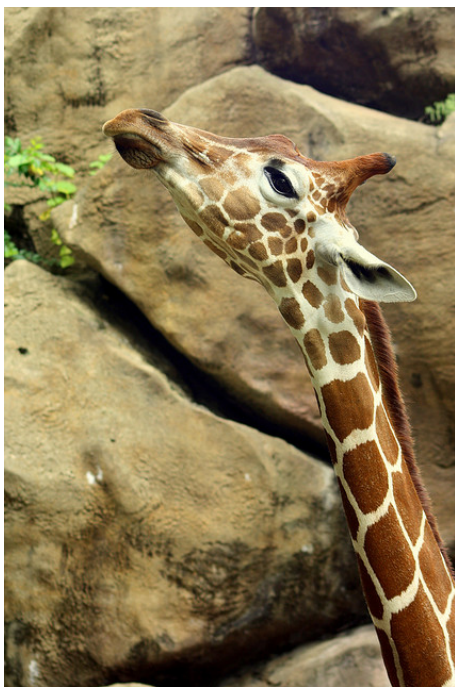


*L*: The food at the restaurant was delicious.

*L* – 1: The little girl and her mother at the food truck.

*L* – 2: a little girl eating a piece of her food at the restaurant.

*L* – 3: The food at the restaurant was delicious.



*L*: The giraffe is a large mammal with a long neck and a spotted coat.

*L* – 1: The giraffe is a species of mammal that is found in africa.

*L* – 2: The giraffe is a species of mammal that is found in the savannas of Africa.

*L* – 3: reciprocal relationship between two animals.

Table 6.6: Examples of captions generated by the models extracting the cls embedding from different layers on images from the MSCOCO dataset.

By observing the results in tables 6.1, 6.2, 6.3 and 6.4 and the examples in tables 6.6 and 6.5 we conclude that changing the layer from where we extract the CLS token gives us different results but is not easy to draw big conclusions, but the main observations are:

- $L$  and  $L - 1$  are the best layers, this can be observed both from the results and the examples but it is not a strong conclusion since for example in table 6.3 we see layer  $L - 3$  with the highest BLEU values but by observing the examples we see that is not advisable to use that layer. The VQA examples in table 6.5 indicate that layer  $L$  is the safest choice.
- The last layer is indubitably the best layer for retrieval, since it gets the best results for both CC3M and MSCOCO.
- Different layers give different pieces of information from the same image. We can observe this from the examples, in the first question of table 6.5 we see that layers  $L$  and  $L - 1$  give rightly attention to man, while the other observe the cars or the ice cream truck.
- For VQA the results in table 6.2 layer  $L - 1$  presents the best results but for a very low margin, but as explained we focus on the examples which indicate that  $L$  is the best choice.

## 6.2 Multiple Vision Tokens

After exploring using visual embeddings extracted from different layers from the ViT we explore how these models behave with a linear transformation leading to multiple embeddings. We are still extracting one visual embedding from the ViT but transforming it with a linear transformation into various embeddings. This follows the process explained in section 3.2.3.1 where to get  $k$  embeddings we get

$$Z_v = h(Y_v) = \text{reshape}(Y_v W + b) \in \mathbb{R}^{k \times 2048} \quad (6.4)$$

where  $W \in \mathbb{R}^{1024 \times (2048k)}$  and  $b \in \mathbb{R}^{1 \times (2048k)}$  obtaining  $Z_v^* \in \mathbb{R}^{1 \times (2048k)}$  before the reshaping operation.

In this section, we explore different values of  $k$  studying how the models behave with different numbers of visual token embeddings. We first experiment for  $Y_v = ViT_{CLS}^L$ , that is extracting the CLS embedding from the last layer and then we experiment for  $Y_v = ViT_{CLS}^{L-1}$ , extracting the CLS embedding from the second to last layer. With this we want to evaluate how the layer influences the model difference in performance for each  $k$ . We analyze for  $k \in \{1, 4, 8, 12, 16, 20\}$ . Although we do present the results for retrieval, we remember that the retrieval is done as explained in section 5.3.1.4 where we have  $W_t \in \mathbb{R}^{d \times q}$  and  $W_i \in \mathbb{R}^{m \times q}$  which transform the [RET] embedding and the image embedding and compares them. This means that the number of tokens  $k$  we are changing

does not directly influence the retrieval but these values differ because of the joint retrieval and captioning training. With this in mind we do report these values but we do not focus on them.

### 6.2.1 Layer $L$

Here we present the results for the different values of  $k$  when we extract the CLS token embedding from layer  $L$ . In figure 6.2 we illustrate the architecture studied in this section.

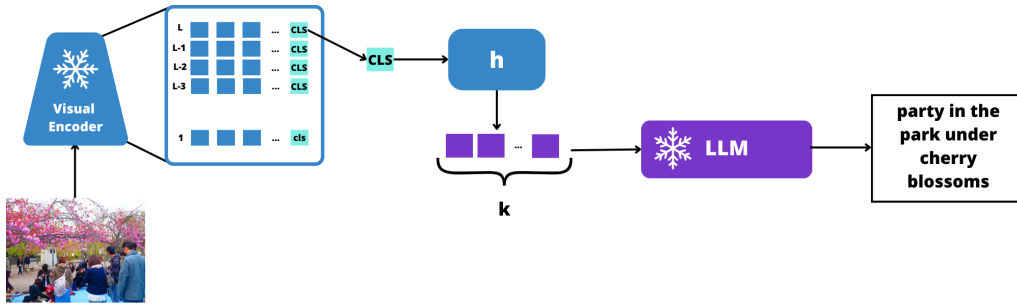


Figure 6.2: Architecture of our model generative part where  $h$  transforms the CLS token of the last layer into  $k$  different tokens, here we change the  $k$ .

Table 6.7: CC3M validation metrics for different numbers of visual tokens with CLS extracted from the last layer

$k$	Captioning					I2T Retrieval		T2I Retrieval	
	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	r1	r5	r1	r5
1	16.30	8.23	4.57	2.74	13.83	23.30	46.02	22.32	44.41
4	23.19	13.54	8.82	6.24	19.40	<b>24.63</b>	47.09	22.63	45.02
8	27.30	16.27	10.80	7.76	22.45	21.65	43.43	19.80	41.01
12	27.83	16.60	11.02	7.90	23.01	23.84	46.45	<b>22.89</b>	45.14
16	27.86	16.78	11.24	8.10	23.40	24.49	<b>47.53</b>	22.85	<b>45.39</b>
20	<b>28.08</b>	<b>17.03</b>	<b>11.46</b>	<b>8.28</b>	<b>23.42</b>	24.05	46.84	22.62	44.57

Table 6.8: VQA validation metrics for different numbers of visual tokens with CLS extracted from the last layer

$k$	overall	yes/no	number	other
1	30.51	60.81	2.82	14.79
4	33.04	63.45	1.39	<b>18.28</b>
8	31.46	61.3	1.99	16.56
12	30.55	59.21	2.29	16.23
16	31.35	61.62	2.62	15.94
20	<b>33.29</b>	<b>65.54</b>	<b>3.02</b>	16.78

Table 6.9: MSCOCO validation metrics for different numbers of tokens extracted from the last layer.

$k$	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr
1	18.87	9.33	4.48	2.22	7.68	21.09	20.64
4	32.59	19.73	11.51	6.75	13.07	33.13	42.12
8	45.24	27.64	16.27	9.70	15.23	35.26	<b>45.76</b>
12	46.81	28.26	16.58	9.85	15.55	34.31	44.15
16	<b>47.76</b>	<b>29.29</b>	<b>17.39</b>	<b>10.30</b>	<b>15.80</b>	<b>35.40</b>	45.74
20	44.44	26.40	15.19	8.93	14.72	33.60	40.81

Table 6.10: MSCOCO retrieval validation metrics for different numbers of token embeddings with CLS token extracted from the last layer.

$k$	I2T		T2I	
	r1	r5	r1	r5
1	14.99	33.97	<b>15.26</b>	34.89
4	15.05	35.43	14.63	34.15
8	13.16	32.12	13.29	31.74
12	<b>15.72</b>	<b>35.89</b>	15.14	34.68
16	15.22	35.11	15.11	<b>35.01</b>
20	8.79	24.13	8.29	22.37

Table 6.11: VQA answers generated by the models with different numbers of tokens extracting the cls embedding from the last layer.




		
<b>Question:</b> What is the man doing in the street?	<b>Question:</b> What are the people in the background doing?	<b>Question:</b> What color is the building?
<b>1 token:</b> he is walking down the street.	<b>1 token:</b> they are watching a skateboarder perform.	<b>1 token:</b> gray.
<b>4 tokens:</b> he's	<b>4 tokens:</b> skateboarding , b: playing video games , c:	<b>4 tokens:</b> grey.
<b>8 tokens:</b> he's walking to work.	<b>8 tokens:</b> skateboarding .	<b>8 tokens:</b> grey.
<b>12 tokens:</b> he is walking.	<b>12 tokens:</b> they are watching the skaters.	<b>12 tokens:</b> it is a light grey colour.
<b>16 tokens:</b> he's selling ice cream.	<b>16 tokens:</b> skateboarding.	<b>16 tokens:</b> grey.
<b>20 tokens:</b> he 's eating ice cream.	<b>20 tokens:</b> skateboarding.	<b>20 tokens:</b> red.

Table 6.12: Examples of captions generated by the models with different number of visual tokens with CLS token extracted from the last layer on images from the MSCOCO dataset.



1 token : The cat is sitting in the car.

4 tokens: a black cat sits in the car.

8 tokens: a cat sits in a car window.

12 tokens: a cat in the car.

16 tokens: cat in the car.

20 tokens: a black cat sitting in a car.



1 token: The food at the restaurant was delicious.

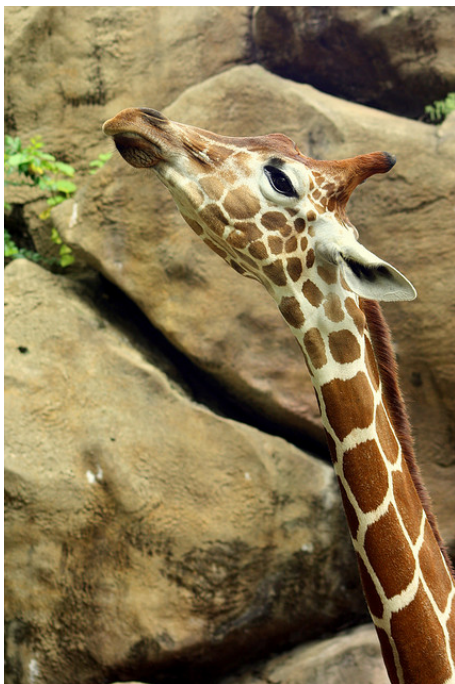
4 tokens: a child eats a hot dog at the restaurant.

8 tokens: person , person , and person enjoying a hot dog at the carnival.

12 tokens: person , a boy , eats a hot dog at the food truck ..

16 tokens: person eating a hot dog at the food tent.

20 tokens: person , eats a hot dog at the food tent.



1 token: The giraffe is a large mammal with a long neck and a spotted coat.

4 tokens: a giraffe in the zoo.

8 tokens: a giraffe in the zoo.

12 tokens: a giraffe in the zoo.

16 tokens: a giraffe in the zoo.

20 tokens: a giraffe in the zoo. photo # 500x333.

We observe that:

- Increasing the number of tokens improves the model's results. This is mainly seen in tables 6.7 and 6.9 where the models with higher number of tokens are the best (20 for CC3M and 16 for MSCOCO).
- 8 tokens is enough. Although we observe improvements with a higher amount of tokens the increase from 1 token to 4 is very significant and from 4 to 8 is still significant while from 8 to more is marginal.
- Models with a high  $k$  (8 or more) may be overfitting. This is what the examples seem to indicate, when we see in table 6.12 captions like " person, person, and person enjoying a hot dog at the carnival." are weird captions that may stem from the poor quality captions from the CC3M training dataset. While the captions for  $k = 4$  like "a child eats a hot dog at the restaurant." seem more natural. But in table 6.11 we see that  $k = 4$  is the one that fails.
- We can conclude that for generative tasks we should choose  $k = 4$  or  $k = 8$  or some value in between that we have not tested.
- For retrieval we can't conclude anything about the number of tokens, which is expected since the number of tokens for retrieval is not changing, therefore the retrieval capabilities of these models are only affected by the captioning-retrieval joint training.

## 6.2.2 Layer $L - 1$

Here we present the results for the different values of  $k$  when we extract the CLS token embedding from layer  $L - 1$ . In figure 6.3 we illustrate the architecture studied in this section.

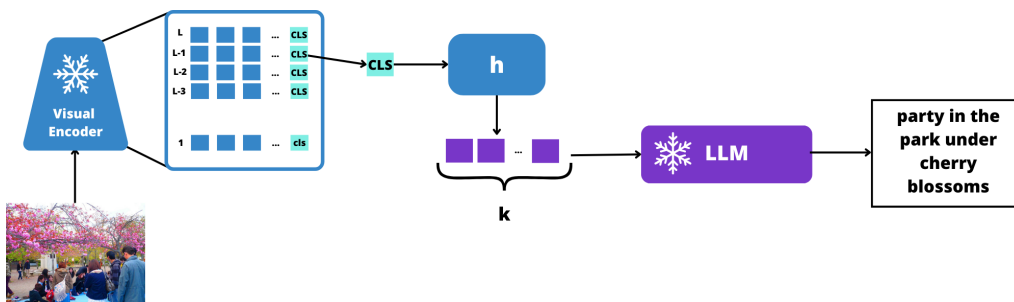


Figure 6.3: Architecture of our model generative part where  $h$  transforms the CLS token from layer  $L - 1$  into  $k$  different tokens, here we change the  $k$ .

Table 6.13: CC3M validation metrics for different numbers of visual tokens with CLS extracted from the second to last layer

$k$	Captioning					I2T Retrieval		T2I Retrieval	
	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	r1	r5	r1	r5
1	16.78	8.45	4.64	2.74	13.86	21.50	43.76	20.53	42.67
4	25.52	14.82	9.62	6.80	20.48	22.02	44.23	20.29	42.34
8	26.78	15.84	10.43	7.43	22.47	21.21	43.14	19.84	41.61
12	27.16	16.16	10.75	7.72	22.49	<b>23.41</b>	<b>45.94</b>	<b>21.43</b>	<b>43.63</b>
16	<b>27.27</b>	<b>16.27</b>	10.81	7.76	23.06	22.10	44.35	21.18	43.42
20	26.96	16.19	<b>10.81</b>	<b>7.79</b>	<b>23.08</b>	21.08	42.58	19.40	40.95

Table 6.14: VQA validation metrics for different numbers of visual tokens with CLS extracted from the second to last layer.

$k$	overall	yes/no	number	other
1	30.53	61.95	1.94	14.20
4	32.02	60.61	2.10	18.19
8	30.70	61.78	1.69	14.74
12	31.71	61.81	2.35	16.57
16	30.96	59.70	2.07	16.75
20	<b>34.36</b>	<b>64.47</b>	<b>3.60</b>	<b>19.60</b>

Table 6.15: MSCOCO validation metrics for different numbers of visual tokens with CLS extracted from the second to last layer.

$k$	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr
1	23.66	12.38	6.20	3.20	9.20	23.12	22.50
4	43.07	25.97	15.06	8.86	14.62	34.46	42.33
8	47.04	28.42	16.53	9.66	15.27	34.55	42.91
12	<b>47.28</b>	<b>28.74</b>	<b>16.83</b>	<b>9.98</b>	<b>15.43</b>	<b>34.76</b>	<b>43.96</b>
16	46.96	28.28	16.46	9.65	15.17	34.43	43.13
20	46.81	28.08	16.25	9.44	15.28	34.00	43.34

Table 6.16: MSCOCO retrieval validation metrics for different numbers of token embeddings with CLS token extracted from the second to last layer.

# of tokens	I2T		T2I	
	r1	r5	r1	r5
1	14.29	33.54	14.22	33.58
4	14.17	32.39	13.39	33.44
8	13.48	32.11	12.89	31.42
12	<b>15.23</b>	<b>34.93</b>	<b>14.94</b>	<b>33.82</b>
16	13.65	32.44	13.66	32.47
20	12.04	29.56	12.34	29.86

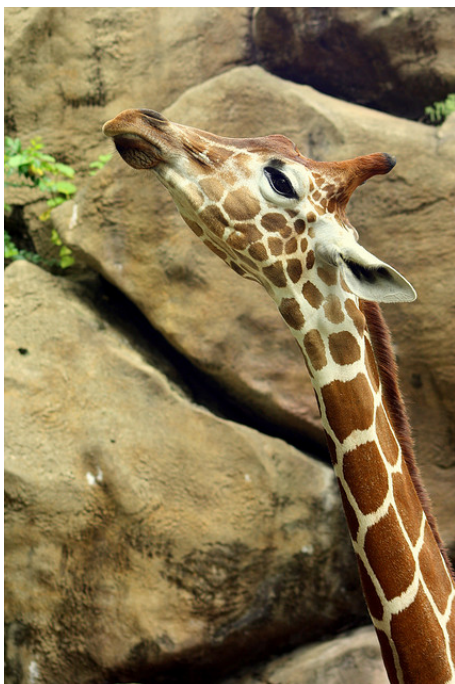
Table 6.17: Examples of captions generated by the models with different number of visual tokens with CLS token extracted from the second to last layer on images from the MSCOCO dataset.



- 1 token : a cat in a car.
- 4 tokens: a cat in a car.
- 8 tokens: a cat sitting in a car .
- 12 tokens: cat in a car.
- 16 tokens: person, the cat, is a bit of a menace.
- 20 tokens: person, a cat in a car.






- 1 token: The little girl and her mother at the food truck.
- 4 tokens: a little girl eating a bite of her mother's food at the festival.
- 8 tokens: person eating a bite of her sandwich.
- 12 tokens: person eating a hot dog at the festival.
- 16 tokens: a little girl enjoying a hot dog at the festival.
- 20 tokens: person eating a hot dog at the festival.



- 1 token: The giraffe is a species of mammal that is found in africa.
- 4 tokens: a giraffe in the zoo.
- 8 tokens: a giraffe in the zoo.
- 12 tokens: a giraffe in the zoo.
- 16 tokens: a giraffe looking at the camera in the zoo. photo # 10248500
- 20 tokens: a giraffe in the zoo.

Table 6.18: VQA answers generated by the models with different numbers of tokens extracting the cls embedding from the second to last layer.

		
<b>Question:</b> What is the man doing in the street?	<b>Question:</b> What are the people in the background doing?	<b>Question:</b> What color is the building?
<b>1 token:</b> he is walking to work.	<b>1 token:</b> the	<b>1 token:</b> gray.
<b>4 tokens:</b> he's trying to get a ride on the streetcar.	<b>4 tokens:</b> skateboarding,	<b>4 tokens:</b> grey, b: green, or c: blue.
<b>8 tokens:</b> he's eating ice cream.	<b>8 tokens:</b> they are watching the skateboarder.	<b>8 tokens:</b> grey.
<b>12 tokens:</b> he 's riding a bike.	<b>12 tokens:</b> they are skateboarding.	<b>12 tokens:</b> grey.
<b>16 tokens:</b> eating ice cream.	<b>16 tokens:</b> skateboarding.	<b>16 tokens:</b> blue.
<b>20 tokens:</b> he is walking to work.	<b>20 tokens:</b> skateboarding.	<b>20 tokens:</b> white.

We observe that:

- Similarly to layer  $L$ , increasing the number of tokens improves the captioning capabilities of the models at least until  $k = 8$ . This time for CC3M (table 6.13) we get the best results for  $k = 16$  and  $k = 20$  while for MSCOCO we get the best results for  $k = 12$ .
- We should choose  $k = 8$ ,  $k = 12$  or some value between that we have not tested. This comes from observing that these values present the best results with a little difference between them and observing the examples both fail in somethings and excel in others. For the captioning examples in table 6.17 they seem both good, the model with 12 tokens correctly understands that the child is eating a hot dog but fails by speculating about a festival that is not an information given by the image. In the VQA examples (table 6.18) the model with 8 tokens seems better since it recognizes that the people are watching and not skateboarding
- For values of  $k$  higher than 1 using layer  $L - 1$  is better than using layer  $L - 1$ , we can observe by comparing tables 6.15 and 6.9 that especially for  $K = 4$  using layer  $L - 1$  is the best option.

- For VQA (table 6.14) the model with  $k = 20$  presents the best results but observing the examples in table 6.18 this seems no to be conclusive.
- For retrieval  $k = 12$  appears to be the best but as explained earlier this is not a strong conclusion since the we always use only 1 token for retrieval.

### 6.3 Other CLS Based Approaches

After analyzing the layers and the number of tokens used we decided to test other methods to using the CLS layers, in these experiments we mixed the earlier ideas questioning whether we could improve these models by using the CLS tokens of more than one layer. We trained the following extra 4 models:

- **L4L** - it stands for last 4 layers, we extract 4 CLS tokens one from each of the last four layers and apply a linear transformation to each, obtaining,

$$Z_v^{L4L} = [LL_0(Y_v^L), LL_1(Y_v^{L-1}), LL_2(Y_v^{L-2}), LL_3(Y_v^{L-3})], \quad (6.5)$$

where  $Y^l = ViT_{CLS}^l(I)$  represents the CLS token embedding from layer  $l$ ,  $LL_i(Y) = YW_i + b_i$ . This is represented in the following figure:

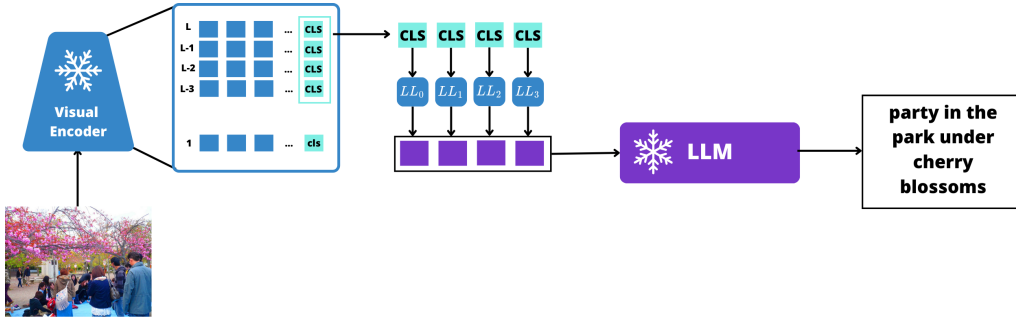


Figure 6.4: Architecture of the generative part of the Last 4 Layers model, where we extract the CLS tokens of the last 4 layers and transform each one with a different linear transformation.

- **L2L** - it stands for last 2 layers, similar to the previous one but only for the last 2 layers, we extract 2 CLS tokens one from each of the last two layers, and apply a linear transformation to each, obtaining,

$$Z_v^{L2L} = [LL_0(Y_v^L), LL_1(Y_v^{L-1})] \quad (6.6)$$

This is represented in the following figure:

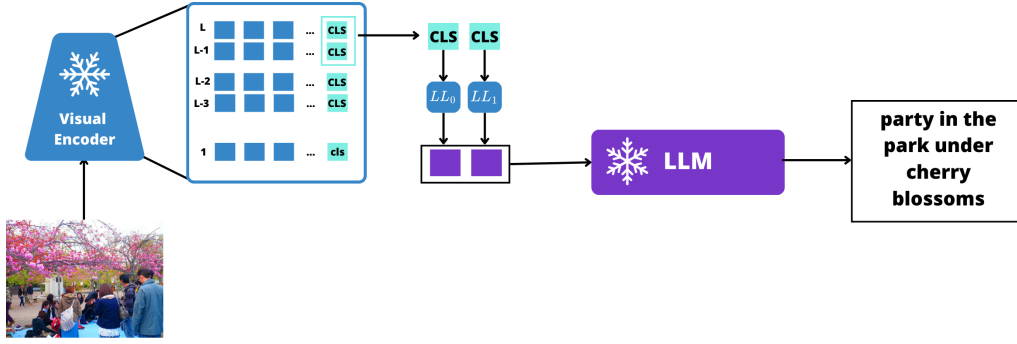


Figure 6.5: Architecture of the generative part of the Last 2 Layers model, where we extract the CLS tokens of the last 2 layers and transform each one with a different linear transformation.

- **L2Lx2** - it stands for last 2 layers times two, it extracts the last two layers but similarly to the methods used in section 6.2 we apply more than one linear transformation to the same vector. In this case we use two linear transformations for each of the two tokens in order put 4 tokens as input of the decoder, since we have concluded that the number of tokens that are input of the decoder is an important variable. With this we obtain,

$$Z_v^{L2Lx2} = [LL_0(Y_v^L), LL_1(Y_v^L), LL_2(Y_v^{L-1}), LL_3(Y_v^{L-1})] \quad (6.7)$$

This is represented in the following figure:

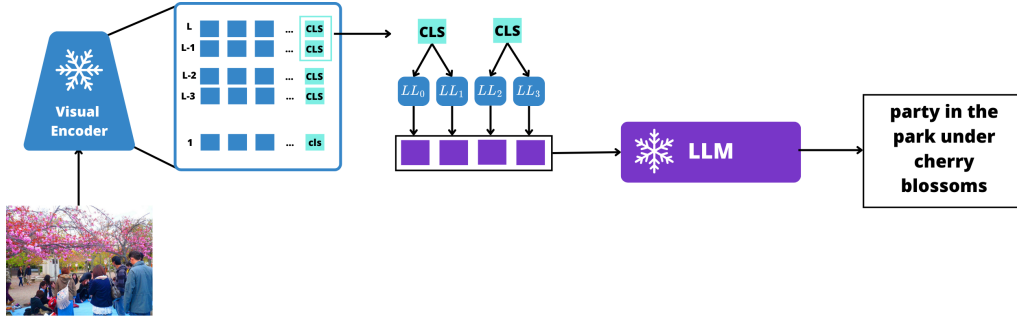


Figure 6.6: Architecture of the generative part of the L2Lx2 model, where we extract the CLS tokens of the last 2 layers and transform each one with two different linear transformations creating 4 embeddings.

- **Comp** - it stands for comparing, we will use "Comp" to identify in the tables but refer to it as the CLS comparing model. This method extracts the CLS tokens from the last two layers and uses them to create two new tokens, one is the difference representing how the CLS tokens changed from one layer to the other and the other is an elementwise multiplication. This should represent the interactions between these two layers, obtaining,

$$Z_v^{Comp} = [LL_0(Y_v^L), LL_1(Y_v^{L-1}), LL_2(Y_v^L - Y_v^{L-1}), LL_3(Y_v^L \odot Y_v^{L-1})], \quad (6.8)$$

where  $\odot$  represents element-wise multiplication. This is represented in the following figure:

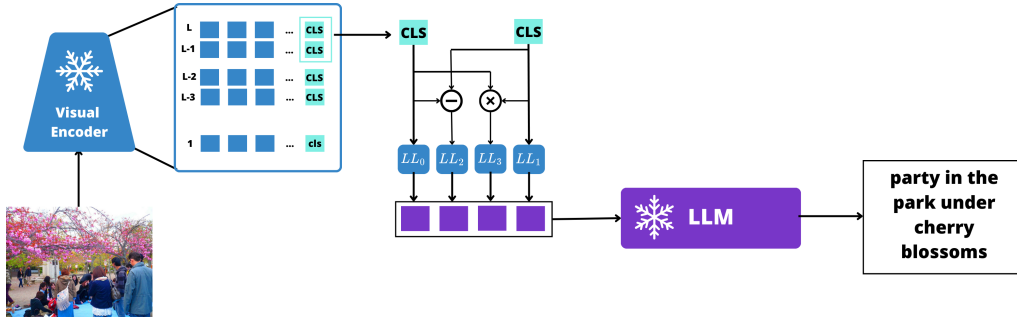


Figure 6.7: Architecture of the generative part of the CLS comparing model, where we use the CLS tokens from the last 2 layers and create 2 new embeddings by subtracting and multiplying element-wise these two. Obtaining 4 tokens where each one is transformed with a different liner transformation.

Since all these models extract more than one CLS token, we use the average of these token embeddings as the retrieval token embedding.

Table 6.19: CC3M validation metrics for different other CLS Based approaches

Method	Captioning					I2T Retrieval		T2I Retrieval	
	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	r1	r5	r1	r5
L4L	24.68	14.33	9.33	<b>6.61</b>	20.12	24.21	47.19	23.05	45.92
L2L	21.76	12.13	7.54	5.12	17.36	24.36	47.10	<b>23.46</b>	45.93
L2Lx2	23.51	13.68	8.95	6.39	18.82	<b>24.49</b>	<b>47.64</b>	23.26	<b>45.99</b>
Comp	<b>25.65</b>	<b>14.61</b>	<b>9.34</b>	6.51	<b>20.73</b>	22.85	45.16	21.11	42.88

Table 6.20: VQA validation metrics for different other CLS Based approaches

Method	overall	yes/no	number	other
L4L	<b>33.31</b>	62.71	2.12	<b>19.22</b>
L2L	32.67	<b>65.04</b>	2.13	16.12
L2Lx2	33.30	63.69	1.94	18.48
Comp	30.79	61.78	<b>2.30</b>	14.75

Table 6.21: MSCOCO validation metrics for different other CLS Based approaches




Method	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr
l4l	45.09	27.56	16.19	9.62	14.95	34.37	<b>43.85</b>
l2l	30.95	18.96	11.06	6.50	12.29	32.27	39.36
l2lx2	40.92	24.91	14.57	8.62	14.14	33.55	42.15
Comp	<b>46.59</b>	<b>28.11</b>	<b>16.35</b>	<b>9.66</b>	<b>15.06</b>	<b>34.43</b>	41.55

Observing the results and examples we conclude that:

Table 6.22: MSCOCO retrieval validation metrics for different other CLS Based approaches

model	I2T		T2I	
	r1	r5	r1	r5
L4L	14.89	35.10	15.08	35.59
L2L	15.25	35.36	<b>15.84</b>	<b>36.08</b>
L2Lx2	<b>15.60</b>	<b>35.89</b>	15.74	35.60
comp	13.87	32.46	13.25	31.68

Table 6.23: VQA answers generated by the models with different numbers of tokens extracting the cls embedding from the second to last layer.

		
<b>Question:</b> What is the man doing in the street?	<b>Question:</b> What are the people in the background doing?	<b>Question:</b> What color is the building?
<b>L4L:</b> he is walking to work.	<b>L4L:</b> they are skateboarding.	<b>L4L:</b> gray, b: green, c: blue, d: black.
<b>L2L:</b> he's trying to get a ride on the streetcar.	<b>L2L:</b> skateboarding,	<b>L2L:</b> grey, b: green, or c: blue.
<b>L2Lx2:</b> he's trying to sell ice cream.	<b>L2Lx2:</b> skateboarder, a person in a hat and a person with a skateboard.	<b>L2Lx2:</b> beige.
<b>Comp:</b> he is walking.	<b>Comp:</b> they are skateboarding.	<b>Comp:</b> gray.

- The number of tokens continues to be important, in the examples we experimented in this section, 3 (L4L, L2Lx2, and Comp) used 4 tokens and 1 (L2L) used 2 tokens, we can see how this leads to L2L being worse than the others and how when we look at tables 6.19 and 6.21 and compare with the previous models, the models with 4 tokens obtain similar results to the other models with 4 tokens and the model with 2 tokens obtains results between the models with 1 token and the models with 4 tokens.
- The Last 4 Layers (L4L) and the CLS Comparing (Comp) models are the ones that perform the best for the generative tasks, both obtaining better results in MSCOCO (table 6.21) than the previous models with  $k = 4$  (only for BLEU and METEOR, not for CIDEr and ROUGE-L, compare with table 6.15).
- If we want only to use 4 tokens we should use Comp for its better quantitative results



L4L: a cat sits in a car window looking out.

L2L: a cat is sitting in a car.

L2Lx2: a cat sits in the car, looking out the window.

Comparison: a black cat sitting in a car.

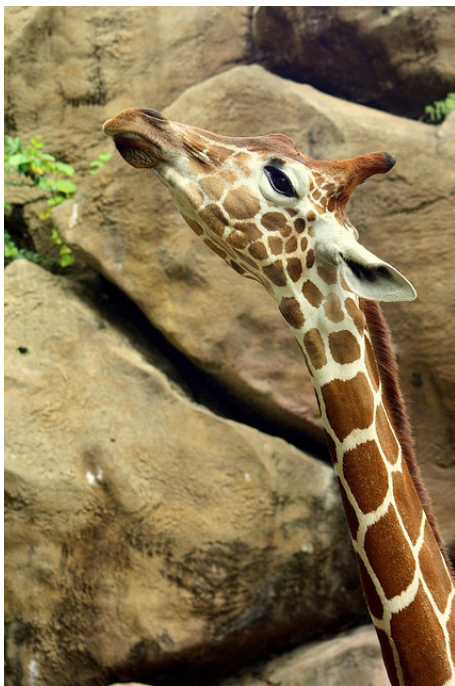


L4L: a little girl and her dad at the carnival .

L2L:

L2Lx2: a little boy takes a bite of his food at the carnival .

Comparison: a child eating a hot dog at the festival.



L4L: a giraffe in the zoo looking up at the sky.

L2L: a giraffe in the zoo.

L2Lx2: a giraffe in the zoo.

Comparison: a giraffe in the zoo.

Table 6.24: Examples of captions generated by the different CLS based approaches on images from the MSCOCO dataset.

but also by observing the examples in tables 6.24 and 6.23 with relatively good captions and answers, even though in table 6.23 it wrongly answers that the people are skateboarding, but that is a recurrent problem in also in the previous models.

- For text-to-image (T2I) retrieval L2L obtains the best results of all the CLS-based models, which may come from combining the CLS token embeddings of the last

two layers.

## 6.4 Summary

From the study presented in this chapter, these are the main conclusions we obtain:

- The last two layers  $L$ , and  $L - 1$  are the best layers to extract the CLS token, depending on the situation one may be better than the other but it is not a drastic change. Extracting from earlier layers is not helpful.
- Independently of the tokens extracted or how they are processed, the most impactful change is the number of tokens that enter the decoder. The other variables can only be compared with the same number of tokens.
- If we want the most from the least amount of tokens  $k$ , we should use  $k = 4$  or  $k = 8$  or some number between that we have not experimented with.
- If we want to use only 4 tokens, using our CLS comparing strategy is probably the best option. The main disadvantage of this approach is that it is not easy to generalize to other numbers of tokens.
- For retrieval tasks the last layer seems to be always preferable, what is expected from the contrastive learning of CLIP.

# ABLATION STUDY: ALL TOKENS BASED MODELS

Now we explore what can be done when we extract all the tokens, which means that  $Y_v \in \mathbb{R}^{257 \times 1024}$ . The most obvious thing to do is to follow what LLaVA [21] does and transform these tokens with a linear transformation, this way, the number of tokens that enter the LLM is the same as the number of tokens on our visual encoder, 257 for our CLIP ViT-L/14 model. However, we were not able to do it as it requires a lot of GPU memory. That way our work is restricted to the use of connecting transformations that reduce the amount of tokens.

## 7.1 Methods

Now, we search for strategies that will reduce the number of tokens effectively. Since we've noted that the number of tokens that enter the LLM makes significant differences in the results, we want to compare the different token reduction methods for similar numbers of tokens so that they can be fairly compared. For this, we chose 4 tokens since the differences from 1 to 4 tokens in previous experiments were the most significant. We experimented with several token reduction approaches to map tokens into a more compact form while retaining essential information. Below, we outline the four main strategies used.

- **Grouped Concatenation with Linear Transformation (Grouped)**

We divided the tokens (excluding the CLS token) into 4 groups and applied a linear transformation on the concatenated embeddings of each group. The transformation maps each group to a vector with the same dimension as the original text embeddings.

Given the embeddings:

$$Y_v = [y_v^0, y_v^1, \dots, y_v^{256}] \in \mathbb{R}^{257 \times 1024} \quad (7.1)$$

We exclude the CLS token  $y_v^0$  and divide the remaining embeddings into 4 groups:

$$Y_v^0 = y_v^1 \oplus \dots \oplus y_v^{64} \in \mathbb{R}^{1 \times 65536} \quad (7.2)$$

$$Y_v^1 = y_v^{65} \oplus \dots \oplus y_v^{128} \in \mathbb{R}^{1 \times 65536} \quad (7.3)$$

$$Y_v^2 = y_v^{129} \oplus \dots \oplus y_v^{192} \in \mathbb{R}^{1 \times 65536} \quad (7.4)$$

$$Y_v^3 = y_v^{193} \oplus \dots \oplus y_v^{256} \in \mathbb{R}^{1 \times 65536} \quad (7.5)$$

Here,  $\oplus$  denotes concatenation. Each group undergoes a linear transformation:

$$LL(Y_v) = Y_v W + b \in \mathbb{R}^{1 \times 2048} \quad (7.6)$$

where  $W \in \mathbb{R}^{65536 \times 2048}$  and  $b \in \mathbb{R}^{1 \times 2048}$ .

Thus, the final embedding is:

$$Z_v^{Grouped} = h(Y_v) = [LL(Y_v^0), LL(Y_v^1), LL(Y_v^2), LL(Y_v^3)] \in \mathbb{R}^{4 \times 2048} \quad (7.7)$$

This is represented in the following figure:

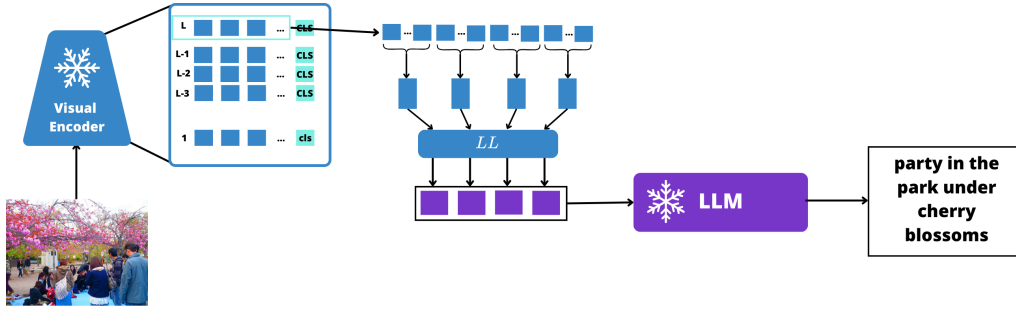


Figure 7.1: Architecture of the generative part of the Grouped Concatenation with Linear Transformation (Grouped) model, where we extract the tokens from layer  $L$  except the CLS and group them in 4 groups where the vectors inside a group are concatenated, obtaining 4 vectors which are transformed by a linear layer.

- **Position-Based Grouping with Linear Transformation (Regions)**

We divided the tokens into 5 groups based on their position in the image. The groups correspond to the 4 quadrants of the image plus an additional central region. Each region contains 64 tokens, with some tokens overlapping between quadrants and the center.

We divide the tokens into:  $Y_v^{TopLeft}$ ,  $Y_v^{TopRight}$ ,  $Y_v^{BotLeft}$ ,  $Y_v^{BotRight}$ ,  $Y_v^{Centre}$

Each group contains 64 tokens. After concatenating the tokens in each region, we apply the same linear transformation  $LL$  as in the previous method:

$$LL(Y_v^{Region}) = Y_v^{Region} W + b \in \mathbb{R}^{1 \times 2048} \quad (7.8)$$

The final representation is:

$$Z_v = h(Y_v) = [LL(Y_v^{TopLeft}), LL(Y_v^{TopRight}), LL(Y_v^{BotLeft}), LL(Y_v^{BotRight}), LL(Y_v^{Centre})] \in \mathbb{R}^{5 \times 2048} \quad (7.9)$$

This is represented in the following figure:

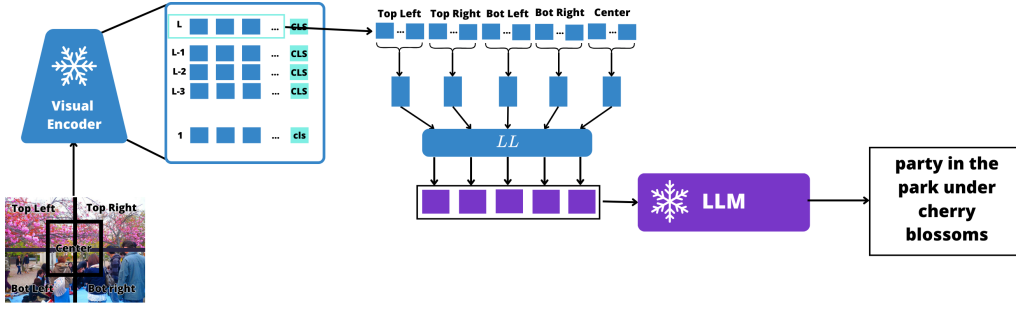


Figure 7.2: Architecture of the generative part of the Position-Based Grouping with linear transformation (Regions) model, where we extract the tokens from layer  $L$  except the CLS and group them in 5 regions corresponding to the region each token represents. The vectors inside a group are concatenated, obtaining 5 vectors which are transformed by a linear layer.

- **Weighted Averaging Instead of Concatenation (W avg)**

Instead of concatenating the tokens, we computed a weighted average of the tokens in each group, which avoids the need to project from a very high dimension (65536).

For each group of tokens:

$$Z_v^0 = LL(W_a[y_v^1, \dots, y_v^{64}]) \quad (7.10)$$

$$Z_v^1 = LL(W_a[y_v^{65}, \dots, y_v^{128}]) \quad (7.11)$$

$$Z_v^2 = LL(W_a[y_v^{129}, \dots, y_v^{192}]) \quad (7.12)$$

$$Z_v^3 = LL(W_a[y_v^{193}, \dots, y_v^{256}]) \quad (7.13)$$

Where  $W_a \in \mathbb{R}^{1 \times 64}$ , and  $[y_v^n, \dots, y_v^{n+63}] \in \mathbb{R}^{64 \times 1024}$ .

The linear transformation  $LL(Y) = YW$  where  $W \in \mathbb{R}^{1024 \times 2048}$  produces:

$$Z_v^{Regions} = [Z_v^0, Z_v^1, Z_v^2, Z_v^3] \in \mathbb{R}^{4 \times 2048} \quad (7.14)$$

This is represented in the following figure:

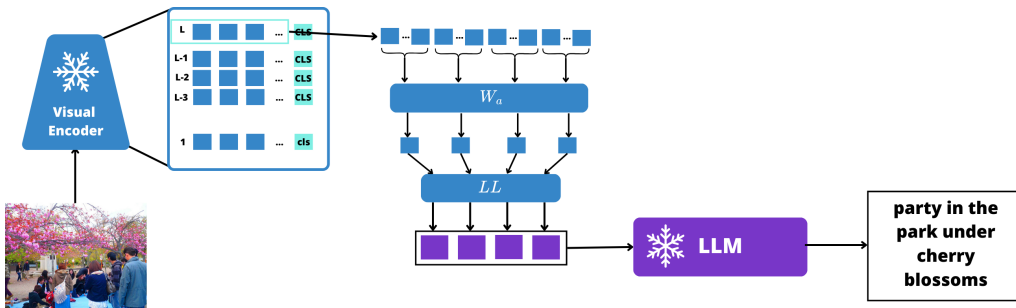


Figure 7.3: Architecture of the generative part of the weighted average (W avg) model, where we extract the tokens from layer  $L$  except the CLS and group them in 4 groups where we obtain a weighted average of the embeddings in each group, obtaining 4 vectors which are transformed by a linear layer.

- **Cross-Attention (CrossAtt)**

Inspired by the Q-Former architecture, we introduced a cross-attention module using 4 query embeddings and 4 attention heads. All image tokens are used as keys and values in this attention mechanism, allowing the queries to focus on different parts of the image.

We use a cross-attention module with query embeddings  $Q \in \mathbb{R}^{4 \times 1024}$ :

$$Z_v^{CrossAtt} = LL_O(head_0 \oplus head_1 \oplus head_2 \oplus head_3) \in \mathbb{R}^{4 \times 2048} \quad (7.15)$$

Each attention head is computed as:

$$head_i = \omega(QK_i^T)V_i \in \mathbb{R}^{4 \times 256} \quad (7.16)$$

where  $\omega$  is the softmax function. The keys and values are computed as:

$$K_i = Y_v W_i^K \in \mathbb{R}^{257 \times 256}, \quad V_i = Y_v W_i^V \in \mathbb{R}^{257 \times 256} \quad (7.17)$$

with  $W_i^K, W_i^V \in \mathbb{R}^{1024 \times 256}$  for each head  $i$ . Finally, the concatenated heads are transformed through  $LL_O(Y) = YW_O + b$  where  $W_O \in \mathbb{R}^{1024 \times 2048}$ . This is represented in the following figure:

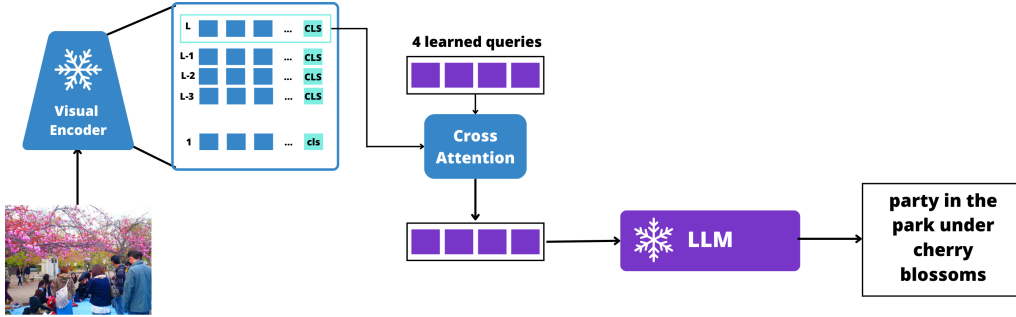


Figure 7.4: Architecture of the generative part of the Cross-Attention model, where we use 4 learned queries which interact with the image tokens from the last layer through cross-attention outputting 4 token embeddings.

Table 7.1: CC3M validation metrics for different methods using all tokens.

Method	Captioning					I2T Retrieval		T2I Retrieval	
	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	r1	r5	r1	r5
CrossAtt	<b>22.83</b>	<b>12.83</b>	<b>8.09</b>	<b>5.59</b>	<b>20.37</b>	<b>20.34</b>	<b>42.36</b>	<b>20.33</b>	<b>41.59</b>
W avg	14.15	7.09	3.83	2.21	12.13	17.23	37.56	15.59	34.84
Regions	15.05	7.70	4.26	2.52	12.24	17.02	36.92	14.69	33.87
Grouped	19.14	9.87	5.55	3.30	14.81	16.27	35.66	13.50	31.90

Table 7.2: VQA validation metrics for different methods using all tokens.

Layer	overall	yes/no	number	other
CrossAtt	<b>30.79</b>	<b>60.45</b>	1.98	<b>15.86</b>
W avg	29.81	58.58	1.64	15.37
Regions	29.28	59.68	0.99	13.64
Grouped	28.58	58.42	<b>3.69</b>	12.45

Table 7.3: MSCOCO captioning validation metrics for different methods using all tokens.

Layer	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr
CrossAtt	<b>46.24</b>	<b>28.53</b>	<b>17.06</b>	<b>10.32</b>	<b>14.62</b>	<b>34.19</b>	<b>43.31</b>
W avg	34.66	20.39	11.62	6.60	10.90	25.83	26.37
Regions	0.73	0.43	0.24	0.13	2.59	5.94	6.41
Grouped	14.88	8.56	4.72	2.60	8.96	28.77	30.34

Table 7.4: MSCOCO retrieval validation metrics for different methods using all tokens.

model	I2T		T2I	
	r1	r5	r1	r5
CrossAtt	<b>14.75</b>	<b>34.39</b>	<b>14.42</b>	<b>34.71</b>
W avg	12.38	30.20	11.40	30.08
Regions	11.24	28.30	11.82	28.98
Grouped	10.74	27.36	10.00	26.64

Table 7.5: Examples of captions generated by different models using various methods to extract all tokens for images from the MSCOCO dataset. The empty captions of the Regions model correspond to images this model returned nothing while trying to caption.



CrossAtt: a cat in a car.

W avg: a cat in a car, looking out the window .

Regions: a cat in a car.

Grouped: a cat is sitting in a car and looking out the window.



CrossAtt: a boy enjoys a meal at the restaurant.

W avg: The boy is eating a hamburger at the restaurant.

Regions:

Grouped: a little girl and her dad at the zoo.






CrossAtt: a giraffe looking up at the sky.

W avg: a giraffe in the zoo.

Regions:

Grouped: a giraffe with a big smile on its face looking at the camera .

Table 7.6: VQA answers generated by the models using different methods to process all visual tokens.

		
<b>Question:</b> What is the man doing in the street?	<b>Question:</b> What are the people in the background doing?	<b>Question:</b> What color is the building?
<b>CrossAtt:</b> he is walking down the street.	<b>CrossAtt:</b> they are watching the skater perform.	<b>CrossAtt:</b> grey.
<b>W avg:</b> he is selling ice cream.	<b>W avg:</b> they are watching a young man skateboarding on a ramp.	<b>W avg:</b> gray , a building with a tower.
<b>Regions:</b> a truck.	<b>Regions:</b> they are watching the crowd of people at the concert.	<b>Regions:</b> the building is painted in a light beige color.
<b>Grouped:</b> he is driving a truck.	<b>Grouped:</b> they are watching the skateboarding competition.	<b>Grouped:</b> the building is a brick building.

Looking at the tables with the results and the examples we observe that:

- The Regions model gives surprisingly low values for captioning in the MSCOCO dataset (table 7.3), this happens because the model is not generating captions for most of the images as can be seen in table 7.5 where it only captioned one of the three images, in total the model only generated 810 captions from the 5000 images, corresponding to about 16% of the images. This model generated text for VQA questions, but don't seem to give accurate answers as we can see in table 7.6, suggesting that the answer comes almost exclusively from the question text and not the image.
- The 4 group model together with the Regions model are the only models that obtain better results for captioning on CC3M than on MSCOCO, this may come from the high dimension (65536) of the concatenated vectors leading to a very large number of parameters being trained that made these models not to be able to generalize to other datasets. For reference the CLS-based models with  $k = 20$  train 42M parameters while the Regions model trains 151M.
- Although it presents the worst values for CC3M the model Weighted average (W avg) obtains reasonable results on MSCOCO (table 7.3), less than the other 4 tokens models but more than the 1 token models, showing that it has good generalizing

capability. We also see in examples from tables 7.5 and 7.6 that its captions and answers although not great are way better than Regions and Grouped.

- The cross-attention model obtains the best results, obtaining better MSCOCO captioning results than most 4-token models and being on par with the CLS comparing model and also generating reasonable captions and answers on the examples of tables 7.5 and 7.6.
- In the different methods used for the all-tokens-based models, Cross Attention is also the best for retrieval. Being always preferable than any of the other methods we tried in this chapter.

From now on when using all tokens we will use the cross-attention module since it obtains the best results for the all tokens based approaches that we tested, results competitive with the CLS only models, and because it is very versatile since we can easily choose the number of tokens we want and we can easily replace it with a more complex module of this kind like the actual Q-Former.

## 7.2 Layer

Here we follow the process done in section 6.1 testing how this model with cross attention as bridging function behaves when we change the layer from which we extract the tokens. For this, we define  $ViT^l$  as the function that extracts all the token embeddings from layer  $l$  of the ViT. We experiment for the last 4 layers, meaning that our function  $f$  will be:

$$Y_v = f(I) = ViT^l(I) \in \mathbb{R}^{257 \times 1024}, l \in \{L-3, L-2, L-1, L\} \quad (7.18)$$

where  $L$  is the last layer. In the previous experiments, we have been using the last layer  $L$ . We illustrate this in figure 7.5.

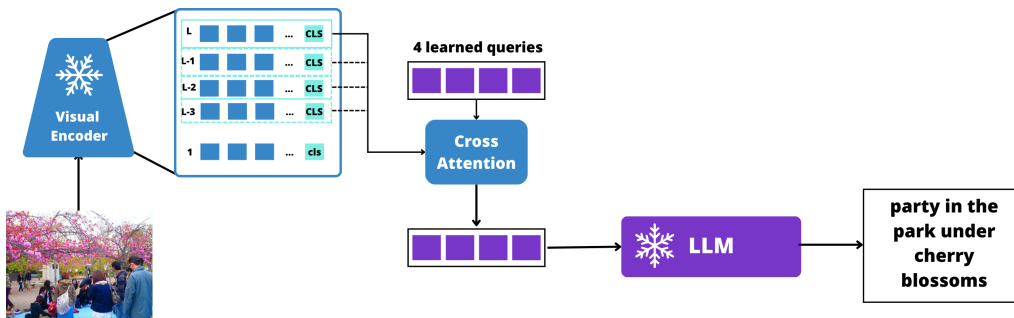


Figure 7.5: Architecture of the generative part of the Cross-Attention model, where we choose from which layer to extract all tokens before they interact with the learned queries through cross-attention. The solid lines show layer  $L$  which is the original layer we used the previous experiments, here we try different layers.

Table 7.7: CC3M validation metrics for different layers using cross-attention.

# of tokens	Captioning					I2T Retrieval		T2I Retrieval	
	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	r1	r5	r1	r5
L	24.07	13.63	8.69	6.06	19.11	<b>23.15</b>	<b>46.13</b>	<b>22.09</b>	<b>44.11</b>
L-1	<b>25.84</b>	<b>14.53</b>	<b>9.18</b>	<b>6.38</b>	<b>20.87</b>	21.00	42.88	20.49	42.17
L-2	24.58	13.88	8.83	6.15	19.81	20.76	43.08	19.94	42.19
L-3	24.24	13.54	8.48	5.81	19.38	20.45	42.75	20.32	42.07

Table 7.8: VQA validation metrics for different layers using cross-attention.

Layer	Yes/no	number	other	overall
L	30.79	15.86	60.45	1.98
L-1	32.32	17.38	62.42	<b>2.30</b>
L-2	<b>32.39</b>	<b>17.40</b>	<b>62.86</b>	1.52
L-3	31.05	15.53	61.77	1.35

Table 7.9: MSCOCO captioning validation metrics for different layers using cross-attention.

Layer	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr
L	<b>46.24</b>	<b>28.53</b>	<b>17.06</b>	<b>10.32</b>	14.62	34.19	<b>43.31</b>
L-1	45.69	26.46	14.90	8.63	14.81	33.00	40.16
L-2	39.24	23.01	12.92	7.33	13.96	32.60	38.75
L-3	46.20	28.03	16.45	9.83	<b>14.99</b>	<b>34.78</b>	43.04

Table 7.10: MSCOCO retrieval validation metrics for different layers using cross-attention.

Layer	I2T		T2I	
	r1	r5	r1	r5
L	<b>14.75</b>	<b>34.39</b>	14.42	34.71
L-1	13.18	31.73	13.10	32.66
L-2	13.78	33.69	<b>15.47</b>	<b>36.08</b>
L-3	13.40	32.75	14.51	34.82



*L*: a cat in a car.

*L* – 1: a cat in a car.

*L* – 2: a cat in a car.

*L* – 3: a cat in a car.

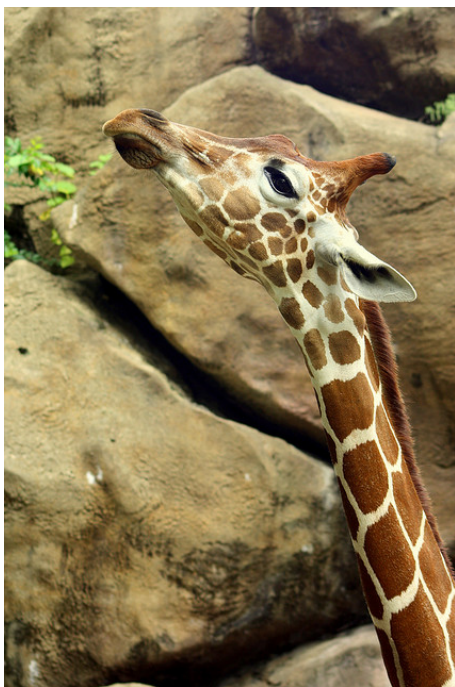


*L*: a boy enjoys a meal at the restaurant.

*L* – 1: a little bit of everything for a little boy.

*L* – 2: a little girl eating a burger at the restaurant.

*L* – 3: a little girl eating a sandwich at a restaurant



*L*: a giraffe looking up at the sky.

*L* – 1: a giraffe in the zoo.

*L* – 2: a giraffe looking at the camera.

*L* – 3: a giraffe looking at the camera.

Table 7.11: Examples of captions generated by the models extracting the cls embedding from different layers on images from the MSCOCO dataset.




		
<b>Question:</b> What is the man doing in the street?	<b>Question:</b> What are the people in the background doing?	<b>Question:</b> What color is the building?
<b>Layer <math>L</math>:</b> he is walking down the street.	<b>Layer <math>L</math>:</b> they are watching the skater perform.	<b>Layer <math>L</math>:</b> grey.
<b>Layer <math>L - 1</math>:</b> he is walking to work.	<b>Layer <math>L - 1</math>:</b> they are watching a young man jump on a skateboard.	<b>Layer <math>L - 1</math>:</b> it's a mix of beige and grey.
<b>Layer <math>L - 2</math>:</b> he is walking.	<b>Layer <math>L - 2</math>:</b> they are watching the skateboarder perform tricks on the ramp.	<b>Layer <math>L - 2</math>:</b> the building is a brick color.
<b>Layer <math>L - 3</math>:</b> he is walking to his car.	<b>Layer <math>L - 3</math>:</b> skateboarding.	<b>Layer <math>L - 3</math>:</b> brick red.

Table 7.12: VQA answers generated by the models with different numbers of tokens extracting the cls embedding from the second to last layer.

Comparing the Cross Attention models extracting all tokens from different layers we observe that:

- There is no clear best layer since all layers perform best in some datasets.
- Giving MSCOCO captioning (table 7.9) the most importance, we see that layer  $L$  is the best choice but not for a significant difference.
- For retrieval layer  $L$  seems the best, although we see that in table 7.10 layer  $L - 2$  performs best for text-to-image retrieval (T2I)

### 7.3 Number of tokens

Finally, we analyze how these models with a cross-attention connection perform if we increase the number of tokens. Increasing the number of tokens in cross-attention means changing the dimensions of our query embeddings  $Q$ . Earlier we explained that  $Q \in \mathbb{R}^{4 \times 1024}$  making  $Z_v \in \mathbb{R}^{4 \times 2048}$ , now when we want  $k$  tokens we have  $Q \in \mathbb{R}^{k \times 1024}$  making  $Z_v \in \mathbb{R}^{k \times 2048}$ . Once again we analyze the performance of these models for different numbers of tokens for tokens extracted from layer  $L$  and from layer  $L - 1$ . As explained earlier, when comparing the number of tokens we do not give a lot of attention to the results on retrieval because they are not directly affected by that. While we do use a

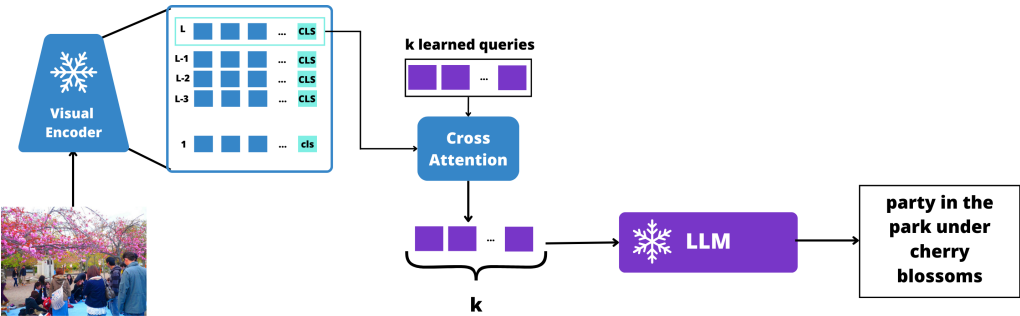


Figure 7.6: Architecture of the generative part of the Cross-Attention model, where we use  $k$  learned queries which interact with the image tokens from the last layer through cross-attention outputting  $k$  token embeddings, here we change the number  $k$ .

cross-attention module to transform embeddings to compare with the embedding of the [RET] token, this uses only one vector as the query.

7.3.1 L

Here we show the results for changing the number of tokens when we extract all tokens from layer  $L$ . We illustrate this in figure 7.6.

Table 7.13: CC3M validation metrics for different numbers of tokens using cross-attention extracting the tokens from the last layer.

$k$	Captioning					I2T Retrieval		T2I Retrieval	
	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	r1	r5	r1	r5
4	24.07	13.63	8.69	6.06	19.11	<b>23.15</b>	<b>46.13</b>	<b>22.09</b>	<b>44.11</b>
8	<b>26.75</b>	<b>15.27</b>	9.79	6.84	21.58	22.40	45.31	21.76	44.08
12	26.37	15.12	9.72	6.78	21.47	20.74	43.25	19.61	41.38
16	26.10	15.02	9.66	6.78	<b>21.73</b>	22.42	45.21	20.97	42.88
20	26.43	15.23	<b>9.85</b>	<b>6.92</b>	21.45	20.22	41.84	19.24	40.54

Table 7.14: VQA validation metrics for different numbers of tokens using cross-attention extracting the tokens from the last layer.

$k$	Yes/no	number	other	overall
4	30.79	15.86	60.45	1.98
8	32.33	17.07	62.66	<b>2.84</b>
12	<b>32.67</b>	<b>17.23</b>	<b>63.77</b>	1.61
16	32.06	16.43	62.86	2.61
20	31.73	16.12	62.55	2.12

Table 7.15: MSCOCO captioning validation metrics for different numbers of tokens using cross-attention extracting the tokens from the last layer.

$k$	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr
4	46.24	28.53	17.06	10.32	14.62	34.19	43.31
8	45.13	26.97	15.58	9.16	15.23	35.01	44.26
12	47.72	28.81	16.87	10.06	15.66	35.13	44.91
16	<b>49.08</b>	<b>29.82</b>	<b>17.51</b>	<b>10.44</b>	<b>15.91</b>	35.41	<b>45.05</b>
20	44.79	27.26	16.09	9.59	15.21	<b>35.93</b>	44.41

Table 7.16: MSCOCO retrieval validation metrics for different numbers of tokens using cross-attention extracting the tokens from the last layer.

$k$	I2T		T2I	
	r1	r5	r1	r5
4	<b>14.75</b>	<b>34.39</b>	14.42	<b>34.71</b>
8	13.85	33.11	14.61	33.93
12	13.53	32.99	13.99	33.75
16	14.51	34.00	<b>14.70</b>	34.27
20	10.69	27.49	10.87	27.30

Table 7.17: Examples of captions generated by the Cross Attention models using different numbers of tokens and extracting all tokens from the last layer on images from the MSCOCO dataset.



4 tokens: a cat in a car.

8 tokens: a cat is sitting in a car.

12 tokens: a black cat sits in the back seat of a car .

16 tokens: a cat is sitting in a car.

20 tokens: a cat in a car.



4 tokens: a boy enjoys a meal at the restaurant.

8 tokens: a little bit of everything for a birthday party.

12 tokens: person , a little girl , eats a hamburger at the restaurant.

16 tokens: a boy eats a hot dog at the food stand.

20 tokens: person , a boy , eats a hamburger at a restaurant.



4 tokens: a giraffe looking up at the sky.




8 tokens: a giraffe in the zoo.

12 tokens: a giraffe in the zoo.

16 tokens: a giraffe in the zoo. the giraffe is a herbivore and the giraffe is the tallest mammal. it is the longest mammal in the world.

20 tokens: a giraffe in the zoo.

Table 7.18: VQA answers generated by the models with different numbers of tokens extracting the cls embedding from the second to last layer.

		
<b>Question:</b> What is the man doing in the street?	<b>Question:</b> What are the people in the background doing?	<b>Question:</b> What color is the building?
<b>4 tokens:</b> he is walking down the street.	<b>4 tokens:</b> they are watching the skater perform.	<b>4 tokens:</b> grey.
<b>8 tokens:</b> he is walking to work.	<b>8 tokens:</b> skateboarders performing a kickflip on a ramp.	<b>8 tokens:</b> grey.
<b>12 tokens:</b> he is waiting for a bus.	<b>12 tokens:</b> they are watching a skateboarder perform a trick.	<b>12 tokens:</b> red.
<b>16 tokens:</b> he 's selling ice cream.	<b>16 tokens:</b> skateboarding. person in the background is skateboarding. person in the background is watching the skateboarder.	<b>16 tokens:</b> the building is white.
<b>20 tokens:</b> he 's walking.	<b>20 tokens:</b> they are watching a skateboarder perform.	<b>20 tokens:</b> it is a building.

Here we report the results for models with cross attention for different values of  $k$  where the tokens are extracted from the last layer ( $L$ ). We observe that:

- The number of tokens continues to be an important variable, but not making a strong difference like in the CLS-based models. We see a not big but significant difference in table 7.13 from  $k = 4$  to  $k = 8$ . But this difference disappears in table 7.15.
- Focusing on table 7.15 we see that using 16 tokens is the best model, obtaining also better results than any CLS-based model. But when looking at the examples in tables, we see that in table 7.17 the captions are good, correctly identifying the hot dog, but we see that in that caption of the giraffe, it has too many not useful words. When looking at table 7.18 we see that the answers are mainly wrong. This suggests that the model is good for captioning but not that good at tasks it was not trained for like VQA.
- The simple 4 token model is a good choice, although it presents worse results than other models in this section, these results are still on par with the best CLS-based

models as mentioned earlier, generating good answers and captions.

- The model with 4 tokens presents the best results for retrieval what once again probably comes from the other models having a lower training loss on captioning, leaving a slightly higher retrieval loss, making it not improve more the retrieval capabilities.

### 7.3.2 L-1

Here we report the results for models with cross attention for different values of  $k$  where the tokens are extracted from the second to last layer ( $L - 1$ ). We illustrate this in the following figure:

Table 7.19: CC3M validation metrics for different numbers of tokens using cross-attention extracting the tokens from the second to last layer.

$k$	Captioning					I2T Retrieval		T2I Retrieval	
	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	r1	r5	r1	r5
4	25.84	14.53	9.18	6.38	20.87	21.00	42.88	20.49	42.17
8	25.81	14.75	9.44	6.59	21.46	<b>21.57</b>	<b>44.13</b>	<b>21.63</b>	<b>43.52</b>
12	<b>26.52</b>	<b>15.26</b>	9.83	6.88	21.71	21.30	43.04	20.87	42.24
16	26.46	15.26	<b>9.85</b>	<b>6.92</b>	<b>21.85</b>	19.03	40.77	18.57	39.69
20	25.05	14.39	9.25	6.48	20.99	18.18	39.27	17.90	38.42

Table 7.20: VQA validation metrics for different numbers of tokens using cross-attention extracting the tokens from the second to last layer.

# of tokens	Yes/no	number	other	overall
4	<b>32.32</b>	<b>17.38</b>	<b>62.42</b>	2.30
8	31.59	16.87	61.25	1.96
12	31.80	17.11	61.20	<b>2.83</b>
16	31.72	17.28	61.00	2.09
20	31.20	16.19	61.13	1.90

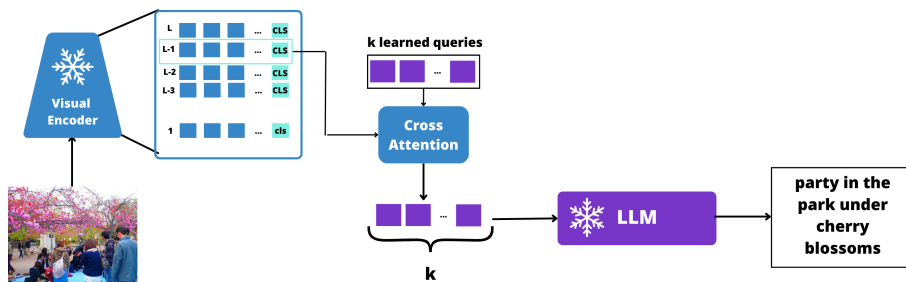


Figure 7.7: Architecture of the generative part of the Cross-Attention model, where we use  $k$  learned queries which interact with the image tokens from layer  $L - 1$  through cross-attention outputting  $k$  token embeddings, here we change the number  $k$ .

Table 7.21: MSCOCO captioning validation metrics for different numbers of tokens using cross-attention extracting the tokens from the second to last layer.

Layer	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr
4	45.69	26.46	14.90	8.63	14.81	33.00	40.16
8	48.64	29.49	17.34	10.33	15.77	35.27	46.31
12	<b>50.02</b>	<b>30.68</b>	<b>18.04</b>	<b>10.81</b>	<b>16.18</b>	36.51	46.96
16	49.36	30.11	17.80	10.58	16.06	<b>36.62</b>	<b>47.07</b>
20	47.88	29.00	17.07	10.20	15.41	34.99	44.31

Table 7.22: MSCOCO retrieval validation metrics for different numbers of tokens using cross-attention extracting the tokens from the second to last layer.

# of tokens	I2T		T2I	
	r1	r5	r1	r5
4	13.18	31.73	13.10	32.66
8	12.97	31.93	<b>14.51</b>	<b>35.24</b>
12	<b>13.91</b>	<b>33.14</b>	14.24	33.78
16	12.51	30.66	12.86	31.88
20	10.92	28.08	11.79	29.45

Table 7.23: Examples of captions generated by the Cross Attention models using different numbers of tokens and extracting all tokens from the second to last layer on images from the MSCOCO dataset.



4 tokens: a cat in a car.

8 tokens: a cat is sitting in a car.

12 tokens: a cat in a car.

16 tokens: a cat in a car, looking out the window.

20 tokens: a cat in a car window.



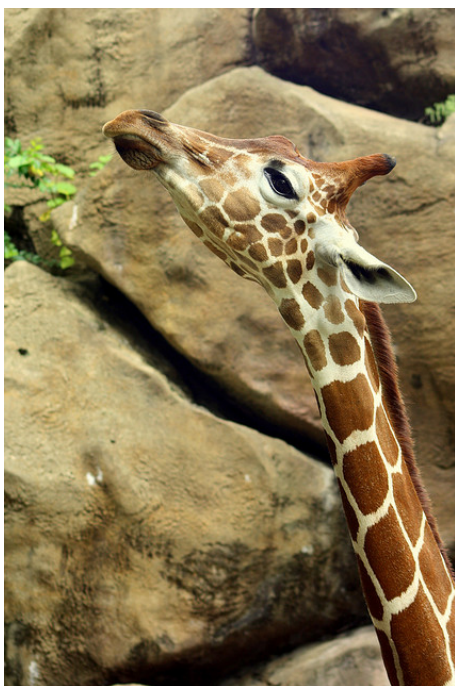
4 tokens: a little bit of everything for a little boy.

8 tokens: person eating a hot dog at the festival.

12 tokens: a little girl eating a hot dog at the carnival.

16 tokens: a little boy and his mother at the food stand.

20 tokens: a little boy enjoying a meal at the restaurant.



4 tokens: a giraffe in the zoo.




8 tokens: a giraffe in the zoo.

12 tokens: a giraffe in the zoo.

16 tokens: a giraffe in the zoo.

20 tokens: a giraffe with a long neck and a long tongue.

Table 7.24: VQA answers generated by the Cross Attention models with different numbers of tokens extracting all tokens from the second to last layer.

		
<b>Question:</b> What is the man doing in the street?	<b>Question:</b> What are the people in the background doing?	<b>Question:</b> What color is the building?
<b>4 tokens:</b> he is walking to work.	<b>4 tokens:</b> they are watching a young man jump on a skateboard.	<b>4 tokens:</b> it's a mix of beige and grey.
<b>8 tokens:</b> he is walking to work.	<b>8 tokens:</b> they are watching the skateboarders.	<b>8 tokens:</b> the building is a mix of red and blue.
<b>12 tokens:</b> he is walking to work.	<b>12 tokens:</b> skaters are doing tricks on the rail.	<b>12 tokens:</b> the building is not painted in any color.
<b>16 tokens:</b> he is walking to work.	<b>16 tokens:</b> they are watching the skateboarders perform.	<b>16 tokens:</b> the building is not a building, it is a museum.
<b>20 tokens:</b> he is selling ice cream to people walking by.	<b>20 tokens:</b> skateboarding.	<b>20 tokens:</b> it is a light blue color.

We observe that:

- When we look at the metrics on CC3M (table 7.19) the number of tokens seems not to make a difference.
- When we look at the metrics on MSCOCO (table 7.21) we see significant differences in the number of tokens. We see that increasing the number of tokens improves the model until  $k = 12$  which is the best model of all the models we experimented.
- The best model is probably between  $k = 12$  and  $K = 16$  since the former has the best values for BLEU and METEOR and the latter the best for ROUGE-1 and CIDEr while presenting similar values in all these metrics. When looking at the examples in tables 7.23 and 7.24 we also see that some captions/answers one performs best and in others the other performs best.
- For retrieval using  $k = 8$  or  $k = 12$  are the best models.

## 7.4 Summary

After the study presented in this chapter and its comparison with the study from the previous chapter, these are the main conclusions:

- Cross Attention is the best of the methods tried both in performance and in versatility since it can take any number of tokens as input and it is easy to increase the amount of output tokens by just increasing the number of query tokens.
- Cross Attention performs better than most CLS-based methods except for the CLS comparing method when using 4 tokens, but has the advantage of easily adding new tokens that the latter does not.
- Similarly to CLS-based models there is not a significant difference between using layers  $L$  or  $L - 1$  but it seems to be better to use layer  $L - 1$  when using more tokens (8 or more) and layer  $L$  when using fewer tokens.
- The best model we tried uses cross-attention with 12 tokens with all tokens extracted from layer  $L - 1$ . This model is also efficient having just slightly more parameters than the CLS model with 4 tokens.

## CONCLUSIONS

### 8.1 Contributions

In this thesis we worked on Generative Vision and Language models, focusing on models that use a frozen image encoder and a frozen text decoder connected by some transformations, following the FROMAGe [14] architecture. In our study and work to advance this kind of Vision and Language Model, our contributions are:

- The development of European Portuguese Vision and Language datasets through machine translation of widely used English datasets. We translated CC3M, MSCOCO, and VisDial, helping train future European Portuguese Vision and Language models.
- Creation of a new European Portuguese Generative Vision and Language Model based on Fromage trained with our translated CC3M dataset.
- Submission of a paper presenting both the translated datasets and the new model.
- An ablation study analyzing different ways to extract and transform visual token embeddings from the ViT encoder to use as input of the LLM decoder, improving the original model.

### 8.2 Take Away Messages

For Portuguese models, we made advancements by translating these datasets and training our new model but it still has some limitations, mainly when we compare with the English models, this probably comes from a lack of larger Portuguese LLM models. The framework we used to train this model allows us to create a better model as soon as there is a new better Portuguese LLM.

From our ablation studies, we conclude that the FROMAGe model can be greatly improved by increasing the number of visual tokens used as input of the LLM, even when just extracting one token from the ViT, showing that there is a high amount of information in the CLS token that cannot be transformed into just one LLM input token and that even

a small number like 4 can bring great improvements. We conclude that cross-attention is one of the best ways to transform a certain amount of tokens into a different determined amount of tokens, being a flexible token reduction strategy that allowed us to train a light Vision and Language model where we can easily choose how many tokens we want to use as input of the LLM.

We could not conclude much about which layer we should extract the token embeddings but it should be either the last layer  $L$  or the second to last layer  $L - 1$ . For retrieval, the last layer is the best, which seems to come from the fact that our image encoder is a CLIP trained with a contrastive learning task. For generative tasks using cross-attention, we concluded that when using multiple tokens we should use the second to last layer.

### 8.3 Future work

Following our ablation studies there are many different changes in the models that would be interesting to study, some examples are:

- Change the ViT - It is important to understand if these models could improve if other image encoders are used, for example by using more recent SOTA models, and how the behaviors and conclusions we observed in these models do or do not generalize for different image encoders. It would also be interesting to compare vision models with different pertaining tasks.
- Change the LLM - Although we already use a different LLM from the one used in the original FROMAGe, it is not certain that LLMs with more parameters would behave the same way these models behaved with Gemma 2B.
- Use more layers - We can add more layers to the transformation using a Multi-Layer Perceptron instead of a linear layer for the CLS-based models or using multiple cross-attentions for all-tokens-based models. It would be interesting to analyze whether this would make a significant change in the models.
- Unfreeze some layers - This is another way to try to improve these models, and see if there is indeed an improvement, by unfreezing, for example, the last layers of the ViT.

For the Portuguese Vision and Language model we could follow the changes tried in these ablation studies to try to improve the model, using more tokens or trying the cross-attention connection.

## BIBLIOGRAPHY

- [1] R. Anil et al. *PaLM 2 Technical Report*. 2023. arXiv: [2305.10403](#) [cs.CL] (cit. on p. 28).
- [2] S. Black et al. *GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow*. Version 1.0. 2021-03 (cit. on pp. 42, 47).
- [3] T. B. Brown et al. *Language Models are Few-Shot Learners*. 2020. arXiv: [2005.14165](#) [cs.CL] (cit. on pp. 28, 41).
- [4] E. Bugliarello et al. *Multimodal Pretraining Unmasked: A Meta-Analysis and a Unified Framework of Vision-and-Language BERTs*. 2021. arXiv: [2011.15124](#) [cs.CL] (cit. on pp. 23–25, 41).
- [5] C. Chambers et al. “FlumeJava: easy, efficient data-parallel pipelines”. In: *Proceedings of the 31st ACM SIGPLAN Conference on Programming Language Design and Implementation*. PLDI ’10. Toronto, Ontario, Canada: Association for Computing Machinery, 2010, pp. 363–375. ISBN: 9781450300193. DOI: [10.1145/1806596.1806638](#). URL: <https://doi.org/10.1145/1806596.1806638> (cit. on p. 37).
- [6] W.-L. Chiang et al. *Vicuna: An Open-Source Chatbot Impressing GPT-4 with 90%\* ChatGPT Quality*. 2023-03. URL: <https://lmsys.org/blog/2023-03-30-vicuna/> (cit. on p. 28).
- [7] H. W. Chung et al. *Scaling Instruction-Finetuned Language Models*. 2022. arXiv: [2210.11416](#) [cs.LG] (cit. on p. 33).
- [8] W. Dai et al. *InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning*. 2023. arXiv: [2305.06500](#) [cs.CV] (cit. on pp. 27, 29, 31, 32, 34).
- [9] A. Das et al. “Visual Dialog”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2017 (cit. on pp. 38, 42, 43, 47, 49).
- [10] A. Dosovitskiy et al. *An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale*. 2021. arXiv: [2010.11929](#) [cs.CV] (cit. on pp. 10, 11, 41).

- 
- [11] D. Jurafsky and J. H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models*. 3rd. Online manuscript released August 20, 2024. 2024. URL: <https://web.stanford.edu/~jurafsky/slp3/> (cit. on p. 9).
- [12] W. Kim, B. Son, and I. Kim. *ViLT: Vision-and-Language Transformer Without Convolution or Region Supervision*. 2021. arXiv: [2102.03334](https://arxiv.org/abs/2102.03334) [stat.ML] (cit. on pp. 25, 41).
- [13] D. P. Kingma and J. Ba. *Adam: A Method for Stochastic Optimization*. 2017. arXiv: [1412.6980](https://arxiv.org/abs/1412.6980) [cs.LG]. URL: <https://arxiv.org/abs/1412.6980> (cit. on p. 47).
- [14] J. Y. Koh, R. Salakhutdinov, and D. Fried. *Grounding Language Models to Images for Multimodal Inputs and Outputs*. 2023. arXiv: [2301.13823](https://arxiv.org/abs/2301.13823) [cs.CL] (cit. on pp. iii, iv, 4, 29–31, 34, 41, 42, 45, 47, 50, 92).
- [15] S. Kudugunta et al. “MADLAD-400: A Multilingual And Document-Level Large Audited Dataset”. In: *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*. Ed. by A. Oh et al. 2023 (cit. on p. 42).
- [16] Y. Lecun et al. “Comparison of learning algorithms for handwritten digit recognition”. In: 1995-01 (cit. on p. 10).
- [17] J. Li et al. *BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation*. 2022. arXiv: [2201.12086](https://arxiv.org/abs/2201.12086) [cs.CV] (cit. on pp. 29, 31).
- [18] T.-Y. Lin et al. *Microsoft COCO: Common Objects in Context*. 2015. arXiv: [1405.0312](https://arxiv.org/abs/1405.0312) [cs.CV]. URL: <https://arxiv.org/abs/1405.0312> (cit. on pp. 36, 37, 42, 43).
- [19] H. Liu et al. *Improved Baselines with Visual Instruction Tuning*. 2024. arXiv: [2310.03744](https://arxiv.org/abs/2310.03744) [cs.CV]. URL: <https://arxiv.org/abs/2310.03744> (cit. on p. 30).
- [20] H. Liu et al. *LLaVA-NeXT: Improved reasoning, OCR, and world knowledge*. 2024-01. URL: <https://llava-vl.github.io/blog/2024-01-30-llava-next/> (cit. on p. 35).
- [21] H. Liu et al. *Visual Instruction Tuning*. 2023. arXiv: [2304.08485](https://arxiv.org/abs/2304.08485) [cs.CV] (cit. on pp. 27, 29, 30, 35, 41, 42, 54, 72).
- [22] R. Lopes, J. Magalhaes, and D. Semedo. “Glória: A Generative and Open Large Language Model for Portuguese”. In: *Proceedings of the 16th International Conference on Computational Processing of Portuguese*. Ed. by P. Gamallo et al. Santiago de Compostela, Galicia/Spain: Association for Computational Linguistics, 2024-03, pp. 441–453 (cit. on pp. 42, 45).
- [23] J. M. Lourenço. *The NOVAtesis L<sup>A</sup>T<sub>E</sub>X Template User’s Manual*. NOVA University Lisbon. 2021. URL: <https://github.com/joaomlourenco/novathesis/raw/main/template.pdf> (cit. on p. i).

- [24] J. Lu et al. “ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks”. In: *Advances in Neural Information Processing Systems*. Ed. by H. Wallach et al. Vol. 32. Curran Associates, Inc., 2019 (cit. on p. 41).
- [25] S. Lu et al. *Ovis: Structural Embedding Alignment for Multimodal Large Language Model*. 2024. arXiv: [2405.20797](https://arxiv.org/abs/2405.20797) [cs.CV]. URL: <https://arxiv.org/abs/2405.20797> (cit. on p. 35).
- [26] OpenAI. *GPT-4 Technical Report*. 2023. arXiv: [2303.08774](https://arxiv.org/abs/2303.08774) [cs.CL] (cit. on p. 28).
- [27] A. Paszke et al. *PyTorch: An Imperative Style, High-Performance Deep Learning Library*. 2019. arXiv: [1912.01703](https://arxiv.org/abs/1912.01703) [cs.LG]. URL: <https://arxiv.org/abs/1912.01703> (cit. on p. 47).
- [28] A. Radford et al. *Learning Transferable Visual Models From Natural Language Supervision*. 2021. arXiv: [2103.00020](https://arxiv.org/abs/2103.00020) [cs.CV] (cit. on pp. 16, 26, 41, 45, 53, 54).
- [29] G. O. dos Santos et al. “CAPIVARA: Cost-Efficient Approach for Improving Multilingual CLIP Performance on Low-Resource Languages”. In: *Proceedings of the 3rd Workshop on Multi-lingual Representation Learning (MRL)*. Ed. by D. Ataman. Singapore: Association for Computational Linguistics, 2023-12, pp. 184–207 (cit. on p. 42).
- [30] R. Santos et al. *Advancing Generative AI for Portuguese with Open Decoder Geroásio PT\**. 2024. arXiv: [2402.18766](https://arxiv.org/abs/2402.18766) [cs.CL] (cit. on pp. 42, 45).
- [31] P. Sharma et al. “Conceptual Captions: A Cleaned, Hypernymed, Image Alt-text Dataset For Automatic Image Captioning”. In: *Proceedings of ACL*. 2018 (cit. on pp. 38, 42, 43).
- [32] H. Tan and M. Bansal. *LXMERT: Learning Cross-Modality Encoder Representations from Transformers*. 2019. arXiv: [1908.07490](https://arxiv.org/abs/1908.07490) [cs.CL] (cit. on p. 26).
- [33] Y. Tang et al. *Multilingual Translation with Extensible Multilingual Pretraining and Finetuning*. 2020. arXiv: [2008.00401](https://arxiv.org/abs/2008.00401) [cs.CL] (cit. on p. 42).
- [34] R. Taori et al. *Stanford Alpaca: An Instruction-following LLaMA model*. [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca). 2023 (cit. on p. 28).
- [35] H. Touvron et al. *Llama 2: Open Foundation and Fine-Tuned Chat Models*. 2023. arXiv: [2307.09288](https://arxiv.org/abs/2307.09288) [cs.CL] (cit. on pp. 41, 42).
- [36] H. Touvron et al. *LLaMA: Open and Efficient Foundation Language Models*. 2023. arXiv: [2302.13971](https://arxiv.org/abs/2302.13971) [cs.CL] (cit. on p. 28).
- [37] A. Vaswani et al. *Attention Is All You Need*. 2023. arXiv: [1706.03762](https://arxiv.org/abs/1706.03762) [cs.CL] (cit. on p. 7).
- [38] P. Wang et al. “OFA: Unifying Architectures, Tasks, and Modalities Through a Simple Sequence-to-Sequence Learning Framework”. In: *International Conference on Machine Learning*. 2022 (cit. on p. 41).

- [39] W. Wang et al. *CogVLM: Visual Expert for Pretrained Language Models*. 2024. arXiv: [2311.03079](#) [cs.CV] (cit. on p. 32).
- [40] Y. Wang et al. *Self-Instruct: Aligning Language Models with Self-Generated Instructions*. 2023. arXiv: [2212.10560](#) [cs.CL] (cit. on p. 33).
- [41] C.-Y. Wu et al. *Sampling Matters in Deep Embedding Learning*. 2018. arXiv: [1706.07567](#) [cs.CV] (cit. on p. 21).
- [42] J. Yu et al. “CoCa: Contrastive Captioners are Image-Text Foundation Models”. In: *Transactions on Machine Learning Research* (2022). ISSN: 2835-8856 (cit. on p. 41).
- [43] A. Zhang et al. *Dive into Deep Learning*. <https://D2L.ai>. Cambridge University Press, 2023 (cit. on pp. 6, 10, 19, 20).
- [44] B. Zhang et al. *Improving Massively Multilingual Neural Machine Translation and Zero-Shot Translation*. 2020. arXiv: [2004.11867](#) [cs.CL] (cit. on p. 42).
- [45] S. Zhang et al. *OPT: Open Pre-trained Transformer Language Models*. 2022. arXiv: [2205.01068](#) [cs.CL] (cit. on p. 41).
- [46] W. X. Zhao et al. *A Survey of Large Language Models*. 2023. arXiv: [2303.18223](#) [cs.CL] (cit. on p. 28).



# 2024 Strategies to Bridge Modalities in Large Vision and Language Models Afonso Simplicio