# Visually Grounded Language Generation: Data, Models and Explanations beyond Descriptive Captions

**Michele Cafagna**

Supervised by Prof Albert Gatt

Co-supervised by Prof Kees van Deemter

Institute of Linguistics and Language Technologies

University of Malta

*A dissertation submitted in partial fulfilment of the requirements for the degree of PhD.*

*Alla famiglia e agli amici*

*In un modo o nell'altro è anche colpa vostra.*

# Acknowledgements

My deepest gratitude goes to my supervisor, Albert Gatt, who has been a great mentor and supporter for both my research and my life in Malta during the last three years. He made my PhD journey a pleasant and rewarding experience. I learned a lot from him about how to conduct, communicate, and appreciate science.

I extend my gratitude to my second supervisor, Kees van Deemter. Despite our in-person interactions being limited to my time in Utrecht, he consistently offered valuable guidance and support.

A special thank-you goes to Lina María Rojas, for supervising my secondment at the Orange Labs, and Marc Tanti for being my reference point at the the University of Malta.

My gratitude also goes to friends and colleagues met in the NL4XAI project. In particular to Ettore Mariotti and Juliette Faille with whom I collaborated during the different secondments, I really enjoyed our discussions. I also thank all the ESRs and the people involved in the project. Without them, my studies and research work would have been way less enjoyable and stimulating.

I also thank all members in the Utrecht Natural Language Processing group and the members of the Institute of Linguistics & Language Technology of the University of Malta.

Finally, I take this opportunity to express my deepest gratitude to my parents, my brother, my grandfather and my closest friends. Their love and support have been important to every twist in my journey.

# Abstract

Vision and Language are two essential capabilities by which we can talk about what we see and communicate it to others, ultimately allowing us to perform tasks, and understand the world. Modeling such interaction is critical to creating agents able to understand, at least to some extent, the world we perceive. This challenge is generally known as multimodal grounding and corresponds to the capability of a model to create meaningful connections between different modalities to solve a task. Ungrounded models do not properly interleave the two modalities yet they can perform well on downstream tasks, leading to misleading and potentially harmful behaviors. Among other fields, Explainable Artificial Intelligence research has moved forward in recent years, proposing methods able to help scrutinize the inner workings of these models and therefore, also assess their grounding capabilities. However, these methods have some relevant limitations, especially on generative models and they are still unpopular in Vision and Language research.

Vision and Language research has mostly focused on performing and evaluating tasks involving the identification and recognition of objects and entities, as they represent the most basic meaningful information represented in a visual scene that can be used as a building block to compose complex multimodal relations, especially on the visual modality. However, in the textual modality, objects represent only a limited amount of linguistic information as language is enriched by words and expressions that do not always correspond to concrete physical objects. Some linguistic expressions can represent complex contexts and situational knowledge that goes beyond the objects visible in the images. For example, describing a picture as a "picnic" (high-level) triggers a whole set of expectations about the scene, making the mention of the objects and entities, totally redundant and uninformative e.g. "people eating food on the grass" (low-level). The latter description is object-centric and it is most likely generated by an automatic captioning system, whereas the former is more human-like and naturally used by humans. The general lack of interest in this relevant aspect by the research community created a potential gap in the overall assessment of the capability of the large-scale models to fully understand the "language", in the "vision and language", preventing a potential gain in terms of overall output quality for multimodal models in generative settings.

In this thesis, we dive into this direction with the aim to discover whether large pre-trained Vision and Language models can handle high-level linguistic descriptions and to what extent they are able to effectively ground them into the visual modality; implications for both language understanding and generation are of interest in this work. Moving away from object-centric descriptions we potentially change the paradigm used to assess multimodal grounding. We analyze potential changes in terms of tasks and evaluation methods introducing an explainability framework designed to complement the currently available tools to assess models' multimodal grounding capabilities in generative settings.

# Contents

# List of Figures

# List of Tables

# List of Abbreviations

# 1

# Introduction

## 1.1 | Vision and Language Modeling

Vision and Language (VL) research lies at the intersection of Computer Vision (CV) and Natural Language Processing (NLP). Its primary goal is to develop algorithms able to jointly understand two intrinsically different modalities, namely textual, and visual. Algorithms of this kind have many applications as they can perform multimodal tasks, like, among others, text-image retrieval, whose goal is to find the most similar image to a given text; or generative tasks like image captioning, where the system automatically generates textual descriptions of images. In order to perform these tasks the algorithms need to meaningfully create connections across the two modalities, this capability is known as *multimodal grounding* and it is essential to handle multimodal tasks.

Given the above, one could state that a well-performing model on a downstream task (e.g. image captioning), has also strong grounding capabilities. However, many studies have shown that VL models can perform well on downstream tasks, though lacking grounding capabilities, by simply exploiting spurious correlations in the data (Ribeiro et al., 2016b). This raises questions regarding current evaluation strategies, which fail to provide an adequate assessment of the grounding capabilities of VL models. In other words, there is a lack of a thorough assessment of the models' capability of aligning linguistic expressions with visual representations, as current evaluation methods are oftentimes task-based. Recent research has been focusing on designing new benchmarks and methods to assess specific linguistic aspects of VL models. This work follows this direction. Our approach to multimodal grounding assessment attempts to provide a different angle, to reason about multimodal grounding, in terms of data, models and tasks. The standard approach to multimodal grounding focuses mainly on the multimodal alignment of visual and textual representations of objects. In this work,

we broaden our focus to the grounding of linguistic expressions in visual representations involving not only objects but also their visual configurations at different levels of abstraction.

# 1.2 | Motivation

## 1.2.1 | Language in isolation is not grounded in the world

When humans use language with communicative intent, such as the intention to convey a message, they rely on concepts and contexts grounded in the world. This capability cannot be obtained by unimodal systems, like a Language Model (LM). As shown in the famous Chinese room experiment (McDermott, 1982; Searle, 1984), an agent that can only access linguistic symbols does not have any experience of the world, as it cannot perceive it. It may be able to produce sensible text and may appear to have "understanding" capabilities. However, it is in fact just relying on linguistic regularities or patterns. Visual data is one of the modalities we can use to model the grounding relationship, enabling models to learn meaning representations that are not entirely unimodal.

## 1.2.2 | High-level expressions are grounded in shared experiences

VL research has always focused on grounding linguistic entities, objects, and attributes into the visual modality (Hodosh et al., 2013a). Textual descriptions of this kind (e.g. "people eating food on the grass") are also called conceptual descriptions, because they focus only on the visible content (Hodosh et al., 2013c), and differ from contextual descriptions (e.g. "people having a picnic") which provide additional situational information, derived from the experience of the world. These kinds of descriptions require cognitive capabilities and are closer to how humans communicate. Successful linguistic communication relies on a shared experience between the interlocutors, which is obtained by building semantic representations on the basis of what we perceive of the world (Bisk et al., 2020).

Perception is the process whereby sensory stimulation is translated into organized experience (Dember et al., 2023). Such a process is influenced by subjective factors that inform what we assume. For example, Figure 1.1[1] shows four subjects: a man, a woman, and two kids; however, an observer would probably describe them as a *family*. They can be perceived both ways, but the latter description is definitely based on a subjective assumption (Van Miltenburg, 2016).

---

[1]Image source `https://pxhere.com/en/photo/1091373`

High-level: *A family.*
Low-level: *A man, a woman, and two kids.*

Figure 1.1: Example of *high-level* and *low-level* descriptions. The latter is an objective description as it describes the only subjects present in the picture, whereas the former conveys additional information regarding the social relationships among the subjects, namely inferring that they are a family.

In other words, we all perceive the same things, but our experiences of them differ according to our beliefs and knowledge of the world, e.g. in Figure 1.1, the visual configuration of the subjects, makes us think of a family.

As a consequence, grounding language into perceptions inform on the knowledge and the biases we use in our interpretation. Such biases can be seen as schemas, put in place to structure and simplify the process of elaborating the perceptual world into experiences and new knowledge. As long as these biases are bound to a shared experience, they are beneficial for communication, as they provide shared patterns and expectations easing the communication that we express through the language.

We constantly rely on these expectations to compose cognitive scripts of scenes and situations helping us in the decision-making process, as described in the well-known script theory (Schank and Abelson, 1975). In language, such expectations are reflected in linguistic expressions that convey additional contextual information. In this work, we shall refer to these as *high-level* descriptions, as opposed to expressions that convey information only regarding the visible content, namely objects, and entities, which we call *low-level* descriptions. Importantly, this distinction varies based on context and the level of abstraction, therefore it is not rigidly binary; instead, it exists and develops along a continuum.

As shown in Figure 1.1, the high-level description provides additional information regarding the social relationships among the subjects. This information is subjective,

as no proof is provided; however, it is a reasonable assumption that can be considered part of common sense. Although, VL datasets are constructed with specific instructions designed to avoid the introduction of such biases, Van Miltenburg (2016) show that they can still be present, confirming that performing inferences is a process that people make routinely.

Current VL approaches aim at grounding object-level, namely low-level expressions, and consider high-level ones a mere source of bias. We argue that object-level information constitutes an essential but limited experience of the world and that high-level information could help VL models enhance their world representation, resulting in a more robust grounding. Nevertheless, the low-level information is critical to make sense of the high-level information; e.g. in Figure 1.1 the reasoning about how the subjects are relating with each other, is triggered by the presence of such individuals and their configuration in the scene.

In this work, we explore this direction, particularly we provide resources and tools to analyze VL models' understanding of high-level information and their impact on the multimodal representation, hoping to foster interest in this line of research.

### 1.2.3 | The role of Explainable AI in assessing VL grounding

Integrating vision and language provides a test-bed for assessing both natural language understanding and goal-oriented visual understanding; indeed VL tasks can demand many disparate CV and NLP tasks to be used simultaneously. Most VL benchmarks capture only a fraction of the requirements of a visual Turing test. In such a test, a system is interrogated on an image with random questions, following a storyline, similar to what humans do when they look at a picture. However, a rigorous evaluation should test each capability required for visual and linguistic understanding independently, so as to help in assessing the right model for the right reason and not for just some spurious correlations. Several studies demonstrated that these systems are affected by dataset biases and a lack of robustness to handle uncommon visual configurations (Choi et al., 2012; Liu et al., 2021a; Thrush et al., 2022; Vedantam et al., 2021). However, although these benchmarks help in identifying such flaws, they do not provide any clue on the reason behind such failures, thus leading to a general lack of transparency in the way grounding is achieved by the models.

In this context Explainable Artificial Intelligence (XAI) is re-emerging as a research trend, with the intent to produce methods able to support, at different granularities, users and researchers in demystifying AI models. However, these days this line of research still places limited emphasis on the VL domain, especially in generative settings.

Moreover, current XAI methods are limited to fine-grained explanations of the linguistic input, namely token-wise, preventing any possible attempt to explain high-level linguistic expressions at the sentence level.

In this work, we focus on a new set of challenges that expand the current view of XAI methods towards generative VL models. with the aim of fostering research in this direction and promoting these methods in the VL field.

## 1.3 | Research Questions

The uncertainty in VL models' performance in different settings mirrors a lack of robustness in their internal representations, due in part to an inconsistent grounding of linguistic expressions in non-visual (perceptual) data and biases in the textual modality towards object-centric descriptions. This thesis focuses on generative vision-to-text models, seeking to thoroughly investigate the following research questions:

1. To what extent do current VL generative models ground high-level information?

2. How can XAI methods be extended to provide a reliable window on model performance with linguistic expressions at different levels?

In order to address research question (RQ) 1 we need to first design a dataset that provides a controlled environment enabling to evaluate VL models against human annotations at multiple levels of linguistic abstraction (high/low-level).

## 1.4 | Roadmap

The remainder of this work is structured as follows. In the next Chapter, we will introduce key concepts related to multimodal grounding and VL models, and how multimodal grounding is enforced. We will give an overview of datasets and tasks and XAI methods. The subsequent three Chapters will be dedicated to answering the two research questions:

- In Chapter 3 will introduce the High-Level (HL) dataset, a new dataset collected to provide high-level descriptions along three axes, namely: scene, actions, and corresponding rationales, aligned with object-centric captions (RQ1).

- Chapter 4 provides a thorough analysis of the capability of VL models to handle high-level descriptions (RQ1).

■ In Chapter 5 we introduce a new framework based on XAI methods designed to assess the grounding capabilities of VL models in generative settings leveraging semantic priors (RQ2).

Finally, we give our conclusions and future work in the last chapter.

## 1.5 | Publications

The following publications are the results of the work conducted during this doctoral study:

1. *What Vision-Language Models "See" when they See Scenes*, **Michele Cafagna**, Kees van Deemter, Albert Gatt, 2021, ArXiv preprint 2109.07301;

2. *Understanding Cross-modal Interactions in V&L Models that Generate Scene Descriptions*, **Michele Cafagna**, Kees van Deemter, Albert Gatt, Proceedings of the Workshop on Unimodal and Multimodal Induction of Linguistic Structures, The 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP2022);

3. *HL Dataset: Visually-grounded Description of Scenes, Actions and Rationales*, **Michele Cafagna**, Kees van Deemter, Albert Gatt, Proceedings of the 16th International Natural Language Generation Conference (INLG2023);

4. *Interpreting Vision and Language Generative Models with Semantic Visual Priors*, **Michele Cafagna**, Lina M. Rojas-Barahona, Kees van Deemter, Albert Gatt, 2023, Frontiers in Artificial Intelligence Journal;

5. *VALSE: A task-independent benchmark for vision and language models centered on linguistic phenomena*, Letitia Parcalabescu, **Michele Cafagna**, Lilitta Muradjan, Anette Frank, Iacer Calixto, Albert Gatt, Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (ACL2022);

6. *TextFocus: Assessing the Faithfulness of Feature Attribution Methods in Natural Language Processing*, Ettore Mariotti, Anna Arias-Duart, **Michele Cafagna**, Albert Gatt, Dario Garcia-Gasulla, Jose Maria Alonso-Moral, IEEE Access Journal;

7. *VILMA: A Zero-Shot Benchmark for Linguistic and Temporal Grounding in Video-Language Models*, Ilker Kesen, Andrea Pedrotti, Mustafa Dogan, **Michele Cafagna**, Emre Can Acikgoz, Letitia Parcalabescu, Iacer Calixto, Anette Frank, Albert Gatt, Aykut Erdem, Erkut Erdem, Proceedings of the Twelfth International Conference on Learning Representations (ICLR2024);

# 1.6 | Code and Resources

**Multimodal Semantic Ablation**   is an ablation method introduced in Cafagna et al. (2021) which allows performing targeted semantic visio-textual ablation. Ablation of text is performed at noun-phrase rather than token level to preserve the input's grammatical correctness. The method detects the noun phrases and removes them from the text generating a new ablated input, generating all the possible combinations of ablated inputs. The visual ablation is performed semantically based on a reference text. The algorithm identifies objects and entities in mentioned in the text and ablates them automatically based on semantic relevancy. The code is available at: `https://github.com/michelecafagna26/vl-ablation`

**HL Dataset**   introduced in Cafagna et al. (2023b) and discussed in this thesis in Chapter 3, is a VL resource aligning object-centric descriptions from an existing VL dataset with high-level descriptions crowdsourced along 3 axes: scenes, actions, and rationales. The HL dataset contains 14997 images from COCO and a total of 134973 crowdsourced captions (3 captions for each axis) aligned with 749984 object-centric captions from COCO. The high-level descriptions capture the human interpretations of the images. These interpretations contain abstract concepts not directly linked to physical objects. Each high-level description is provided with a confidence score, crowd-sourced by an independent worker measuring the extent to which the high-level description is likely given the corresponding image, question, and caption. The higher the score, the more the high-level caption can is close to commonsense (on a Likert scale from 1-5). The dataset is officially released at: `https://github.com/michelecafagna26/HL-dataset/tree/main;` and available on the HuggingFace Hub at: `https://huggingface.co/datasets/michelecafagna26/hl`.

A **dataset explorer** for the HL Dataset is provided at: `https://huggingface.co/spaces/michelecafagna26/High-Level-Dataset-explorer`.

The dataset is further extended with the **HL-Narratives**; a new dataset containing synthetic narratives of the image, generated by combining the three axes of the HL Dataset. The dataset is generated by an automatic hybrid procedure involving an LLm and human feedback. The **HL-Narratives** is available at: `https://huggingface.co/datasets/michelecafagna26/hl-narratives`.

We release a total of 14 **baseline models** for both datasets, available at: `https://huggingface.co/michelecafagna26`.

**VL-SHAP**   is an explainability method for VL generative models introduced in Cafagna et al. (2023a) and thoroughly discussed in Chapter 5- The method is based on the KernelSHAP method and implements several features: (1) a deterministic approximation method to compute Shappley values; (2) sentence-based visual explanations allowing the explanation of the whole generated caption rather than token-based explanations, (3) explanations exploiting visual semantic priors learned by the model (4) The method is model-agnostic and efficient. The code is available at: `https://github.com/michelecafagna26/vl-shap`.

# Background

## 2.1 | What is Multimodal Grounding?

**What does grounding mean?** The problem of making the semantic interpretation of a symbol intrinsic to a system is known in cognitive science as *the symbol grounding problem* (Harnad, 1990). Harnad (1990) in his influential work poses a fundamental question: *How can the meaning of symbols, manipulated solely based on their arbitrary shape be grounded in anything other than symbols?*

This problem is analogous to Searle (1984)'s thought experiment, consisting of trying to learn Chinese only from a Chinese dictionary. It is hypothetically possible to become fluent in Chinese by finding patterns in the symbols, but how could you understand the meaning of those symbols without connecting them to the world?

A plausible solution to the symbol grounding problem is to rely on non-symbolic representations, namely representations built upon non-arbitrary structures. Representations of such kind, referred to as "iconic representations" by Harnad (1990), originate from the elaboration of the sensory signal perceived from objects and events of the world; thus they are good candidates to connect symbols to nonsymbolic world representations. Therefore, a system able to perform grounded symbol manipulations would be driven not just by symbolic tokens, e.g. linguistic tokens, but by non-arbitrary iconic representations, e.g. visual inputs, in which those symbols are grounded (see also Bender et al., 2021).

**What does multimodal mean?** With the term modality, we usually refer to some form of representation conveying information of the world through specific channels, such as visual, auditory, and textual. A multimodal system must process and relate several

Figure 2.1: Image from the MS-COCO 2014 validation set. One reference caption is: *a man in a chefs hat chopping food* [reproduced verbatim from the dataset].

modalities together. This requires acquiring a semantic representation of concepts related to the world, based on such modalities.

Multimodal grounding can be defined as the capability of a multimodal system to find meaningful connections across multiple modalities.

In this work, we focus on vision and language models (VL), namely systems processing images and text. Grounding, in this context, means building meaningful links between linguistic expressions and visual features. To achieve this goal, VL models process the textual and visual modalities in modules called *encoders* which create numerical representations of the inputs. Encoders will be fully described in Section 2.2.2. Grounding is performed by enforcing multimodal interactions between these numerical representations. This happens mainly through two mechanisms: *attention* and *pre-training objectives*. The former is an architectural component of the model and can be implemented in several variations. These will be thoroughly discussed in Section 2.2.3. The latter is a loss function optimized by the learning algorithm. The most relevant pre-training objectives will be discussed in depth in Section 2.2.5.

Grounding language in visual representations is fundamental for human communication, as such representations cannot be acquired uniquely by linguistic symbols (Beinborn et al., 2018). This process is not trivial as vision and text have different characteristics, which also implies that one modality can provide complementary, redundant, or conflicting information with respect to the other. For example Figure 2.1 shows a bird on the shoulder of the chef which is not mentioned in the corresponding caption. The captions may have different levels of abstractions and granularity e.g. a description such as "a man cooking" for Figure 2.1 though quite abstract, provides enough context to describe the image.

10

Multimodal grounding is critical for ensuring trust and transparency in VL systems. While downstream task performance is commonly used for evaluation, true grounding, essential for understanding visual content, cannot be guaranteed solely based on task success. As Ribeiro et al. (2016b) showed in their influential work, a model performing well on a specific task does not necessarily have grounding capabilities, as it may solve the task purely based on spurious correlations in the data. Neglecting proper grounding assessment can lead to misleading and potentially harmful behaviors, especially in domains such as healthcare and navigation. Therefore, performing a rigorous assessment of the multimodal grounding is fundamental to ensure trustworthiness and transparency in VL systems.

Moreover, current VL models can be pre-trained on large amounts of data and easily adapted to many different tasks and domains. This incredible flexibility of pre-trained VL models raises questions regarding the effective capability of these models to ground and generalise representations of unknown data, making the proper assessment of multimodal grounding a critical and necessary requirement.

### 2.1.1 | Outline of this Chapter

In the remainder of this Chapter, we will set the stage for the whole work. We will provide the background on key aspects related to VL modeling and grounding helpful to understand the research problems tackled in this work. However, a related work section will be included with each Chapter to discuss the relevant literature there.

In Section 2.2, we will give an overview of current multimodal architectures and the strategies implemented to enforce grounding across modalities. Section 2.3 will provide a brief overview of the main VL tasks and datasets, and evaluation techniques. The end of this Section is dedicated to introducing the main concepts related to XAI methods, emphasizing the potential benefits of these methods in assessing the extent to which multimodal grounding is successful.

## 2.2 | Vision and Language Modeling

VL models are neural networks that can process both visual and textual inputs such as images and captions. In this Section, we will discuss the main VL architectures proposed for a wide range of VL tasks such as image-sentence matching, image captioning, visual question answering, and visual dialog. We review the main datasets, tasks and strategies implemented to enforce multimodal grounding, highlighting their strengths and limitations.

This Section is structured as follows: in Section 2.2.1 we discuss macro-architectures, then we go into detail in Section 2.2.2 on the major architectural building blocks while Section 2.2.4 provides an overview of the major micro-architectural differences. In Section 2.2.3 and Section 2.2.5 we discuss respectively modality fusion mechanisms and pre-training objectives. In light of what was discussed in the preceding sections, Section 2.2.6 briefly summarises the major strategy implemented to enforce multimodal grounding in VL modeling. Section 2.2.7 concludes with a brief historical overview marking the transition from task-specific to large-scale models.

## 2.2.1 | Architectures

VL architectures can be grouped into two categories, namely single- and dual-stream models (Long et al., 2022). **Single-stream** models fuse the two modalities at an early stage (e.g. UNITER (Chen et al., 2020b), VisualBERT Li et al. (2020c)) to produce multimodal representations. **Dual-stream** models (e.g. CLIP Radford et al. (2021b), ViLBERT (Lu et al., 2019)), process the modalities separately, and then they enable multimodal interactions only at a later stage. This allows the creation of intermediate unimodal representations that can be fused (e.g. LXMert (Tan and Bansal, 2019), CoCa (Yu et al., 2022)) or aligned (e,g. ALIGN (Jia et al., 2021), Florence (Yuan et al., 2021)) together to solve multimodal tasks.

In terms of performance, they have proven to perform roughly on par (Bugliarello et al., 2021a). However, these two architectures feature relevant pros and cons. The main advantage of single-stream models is their efficiency as they have usually a smaller number of parameters with respect to the dual-stream counterpart. However, the dual-stream models are more versatile, as their unimodal encoders can be used separately to perform unimodal tasks, or combined to compose more complex architectures, as also shown by Singh et al. (2022). This great reusability may balance the higher resource cost linked to running larger models.

## 2.2.2 | Encoders

**Textual Encoders**   VL models encode the textual input using Transformer-based encoders (Vaswani et al., 2017b) following BERT (Devlin et al., 2019a), ROBERTA (Liu et al., 2019) and GPT (Brown et al., 2020a). The latter differs from the first two both in that it has an autoregressive self-attention mechanism (bidirectional for BERT and ROBERTA) and in its pre-training objectives. Nevertheless, they all share the same general architecture.

The input sentence is first split into sub-words (Kudo and Richardson, 2018) and then bounded by special tokens at the beginning and the end of the sentence. Occasionally, another special separator token is added to separate multiple sentences in the same sequence. The text is encoded into position-aware learnable embeddings that are then fed into successive text-specific layers before fusion (dual-stream) or directly into a fusion module (single-stream) in order to obtain multimodal representations. Pre-trained BERT-based textual encoders produce semantic textual representation and can provide a good starting point for pre-training on multimodal tasks. In some models, these representations are used with minimal or no fine-tuning, as in Flamingo (Alayrac et al., 2022) where a large pre-trained LM is kept frozen during multimodal few-shot tuning, achieving impressive results on downstream tasks.

**Visual Encoders**   process the visual input to produce meaningful features encoding the relevant visual information. Gan et al. (2022) identifies three types of vision encoders:

- **Object Detector (OD)** models are trained on CV datasets such as Objects365 (Shao et al., 2019) and Visual Genome (VG) (Krishna et al., 2017a) and OpenImages (Kuznetsova et al., 2020) on the object detection task. In VL they are used to provide pre-trained visual features as they provide meaningful representations for objects, entities, and their spatial location in the image, as in Bottom-up and top-down attention (BUTD) (Anderson et al., 2018). The FasterRCNN (Ren et al., 2015b) is the most used OD; it generates features representing bounding boxes of objects detected in the image. Other popular OD encoders are based on the ResNet's architecture (He et al., 2015a) such as VinVL (Zhang et al., 2021), where the use of a more capable visual backbone on CV tasks, was shown to have a positive impact on the performance of VL tasks.

- Plain **CNNs** are usually ResNets pre-trained on the image classification task, on datasets such as ImageNet (Deng et al., 2009). Differently from ODs, they produce visual representations linked to semantic concepts (Rombach et al., 2020), capturing visual features (Collins et al., 2018) that can be extracted from the intermediate layers. Despite the lower popularity with respect to the ODs, they have comparable performance to OD-based visual backbones. Examples of VL models adopting CNN-based visual backbones are PixelBERT (Huang et al., 2020), SimVLM (Wang et al., 2021) and the ResNet version of CLIP (Radford et al., 2021b).

- **Vision Transformer (ViT)** (Dosovitskiy et al., 2020b) is a Transformer model adapted to process visual inputs. Since its introduction, it has shown comparable perfor-

mance to CNN-based models, achieving state-of-the-art performance in CV and VL tasks with the benefit of being more efficient to train. In a ViT, the image is split into equally sized square patches, which are flattened out. Then they are linearly embedded with position embeddings and fed into the vanilla multilayer Transformer. Since its introduction ViT-based visual backbones have gained popularity in VL, among the many models we mention DeiT (Touvron et al., 2022), BEiT (Bao et al., 2021a), Swin Transformer (Liu et al., 2021b) and CLIP-ViT (Radford et al., 2021b).

## 2.2.3 | Multimodal Fusion

The success of Transformer models is due to the attention mechanism. The attention allows to capture short- as well as long-distance relations in the input sequence. The idea behind the attention is that all the elements in a sequence should interact with each other to some extent, though this varies depending on the training objective. For example, in the autoregressive scenario, more common in VL generative models, each input element is forced into attending only the preceding ones. This interaction is regulated by learned parameters that dynamically modulate the extent to which an input feature affects another one, namely the *attention weight*. This mechanism is naturally suitable for performing modality fusion and enforcing multimodal grounding, as it allows inter-modal interactions between the modality-specific features. In this Section 2.2.3.1 we will provide a formal definition of the general attention mechanism implemented in Transformer layers, usually referred to as *self-attention*. In Sections 2.2.3.2 and 2.2.3.2 we will discuss how this mechanism is implemented and adapted to different VL architectures.

### 2.2.3.1 | Transformer Layer

A Transformer-based architecture is composed of a stack of Transformer layers. Following the Bugliarello et al. (2021a) notation, a Transformer layer can be decomposed in two main components: a Multi Head Attention Block (MAB) and a Feed Forward Block (FFB):

**Attention Head**   Given $N_q$ query vectors, each of dimension $d_q$, $Q \in \mathbb{R}^{N_q \times d_q}$, and $N_v$ key-value pairs $K \in \mathbb{R}^{N_v \times d_q}$, $V \in \mathbb{R}^{N_v \times d_v}$, the attention head $Att(Q, K, V)$ maps queries to output vectors through a scaled dot product:

$$Att(Q, K, V) = \omega(QK^T) \tag{2.1}$$

where $\omega$ is a row-wise scaled softmax: $\omega_i(\cdot) = softmax(\cdot/\sqrt{d_q})$. The attention head computes a similarity matrix between queries and keys: $S = QK^T \in \mathbb{R}^{N_q \times N_v}$, then weighted with $V$.

**Multi-head Attention**  projects $Q, K, V$ into $H$ different matrices and then computes the attention (Eq. 2.1) on each projection. The $H$ different outputs are then concatenated togeter ($\|$) and projected with another transformation $W^O$:

$$MHA(Q, K, V) = [O_1\|...\|O_H]W^O,$$
$$\text{where}\quad O_h = Att(QW_h^Q, KW_h^K, VW_h^V) \tag{2.2}$$

The projection matrices $W_h^Q, W_h^K, W_h^V$ and $W^O$ are learned parameters. Note that usually, $d_q = d_v = d$ and $da = d/H$, therefore $W^O \in \mathbb{R}d \times d$ and $W_h^Q, W_h^K, W_h^V \in \mathbb{R}^{Nxd}$.

**Multi-head Attention Block**  can now be defined as:

$$MAB(X, Y) = LN(X + MHA(X, Y, Y)) \tag{2.3}$$

where $X, Y \in \mathbb{R}^{Nxd}$ are given inputs and $LN$ is layer normalization (Ba et al., 2016),. with $X$ being the residual connection.

**Feed forward Block**  for an input matrix $M \in \mathbb{R}^{N \times d}$ is given by:

$$FFB(M) = LN(M + ReLU(MW_1)W_2) \tag{2.4}$$

Where the projection matrices $W_1, W_2^T \in \mathbb{R}^{d \times d_f}$ are learnable parameters and $ReLU$ is the Rectified Linear Units activation function (Agarap, 2018).

**Transformer Layer**  performing *self-attention* can finally be defined as:

$$f_\theta = FFB(MAB(X, X)) \tag{2.5}$$

A **Transformer Encoder** is built by stacking up a sequence of Transformer layers namely:

$$E(X) = f_{\theta_L} \circ ... \circ f_{\theta_1}(X) \tag{2.6}$$

### 2.2.3.2 | Attention-based Fusion

In VL modeling Transformer-based models build upon two main implementations:

15

Figure 2.2: Visualization of attention mechanism implemented in Transformers-based models. 2.2a single-stream using self-attention; 2.2c dual-stream using co-attention 2.2b dual-stream using self-attention. The modalities are processed separately and fused at a later stage with a different mechanism, e.g. a pre-training objective. Figure elaborated from Bugliarello et al. (2020).

**Self-Attention Transformer Layer**   Self-Attention is the traditional attention mechanism employed in the Transformers which the single-stream architectures rely on. The visual and the textual features are concatenated and processed together in the Transformer block. Therefore, the input matrix $X_{VL}$ is defined as follows: $X_{VL} = [X_V||X_L] \in \mathbb{R}^{Nxd}$ with $N = N_V + N_L$, where $V$ and $L$ denote visual and textual respectively. As shown in Figure 2.2a, a single-stream model executes Eq. 2.6 with no modification, namely $E(X_{VL})$

In the first layer, the modalities start interacting, producing multimodal representations at each intermediate layer. Models implementing self-attention include UNITER (Chen et al., 2020b), VisualBERT (Li et al., 2020c), LEMON (Hu et al., 2022b), GIT (Wang et al., 2022a) and Unicoder-VL (Li et al., 2020a).

**Co-attention Transformer Layer**   Co-attention is widely used in dual-stream architectures. The visual and textual inputs are processed into separate Transformer blocks, and then a cross-attention technique is used to force multimodal interactions. It allows unimodal representations to affect each other at multiple stages. Assume $X_V \in \mathbb{R}^{N_V xd}$ and $X_L \mathbb{R}^{N_L xd}$ the visual and the textual input respectively; the inputs are processed in separate encoders $E_V(X_V), E_L(X_L)$. To perform co-attention, the $MAB$ computes the similarity matrix $S$ over queries and key-value pairs of different modalities. Without loss of generality, let's consider the visual encoder $E_V$. Following Eq. 2.3, we obtain $MAB_V(M_V, M_L)$. Hence, Eq. 2.2 is computed as:

$$MHA_V(Q_V, K_L, V_L) \tag{2.7}$$

where $Q_V$ are the queries computed over the visual input and $K_T$, $V_T$ are the key-value pairs computed over the textual signal. The same applies to the textual encoder $E_L$. A schematic representation of the method is shown in Figure 2.2c.

This mechanism is used in models like VilBERT (Lu et al., 2019), LXMert (Tan and Bansal, 2019), (Yu et al., 2022), BLIP (Li et al., 2022c) and ALBEF (Li et al., 2021b) to name a few.

In some dual stream models like CLIP (Radford et al., 2021b) and ALIGN (Jia et al., 2021), the encoders use traditional self-attention Transformer layers, keeping the computation of the two modalities separate in the whole architecture. A schematic example is shown in Figure 2.2b. In these models, the fusion of the two modalities is performed by a multimodal pre-training objective, namely the dot-product between two modality-specific vector representations. This mechanism forces the unimodal encoders to project their representation into a joint multimodal space. Differently from the attention-based fusion mechanism, in this method, the multimodal grounding enforcement relies entirely on the optimization of a multimodal objective task. The pre-training objective and their role in grounding enforcement will be tackled more in depth in Section 2.2.5.

## 2.2.4 | Encoder-Only vs. Encoder-Decoder

A large number of VL models follow an encoder-only architecture (Figure 2.3a), where the two modalities, are fused in a fusion model that then generates the final outputs. This design is naturally suitable for VL understanding tasks like VQA and Visual Spatial Reasoning (VSR) (Liu et al., 2023). However, it can be adapted to generative vision-to-text tasks, by manipulating the attention masks to autoregressively generate sequences of text, resulting in a *Decoder-only* architecture (Figure 2.3b), such as in Oscar (Li et al., 2020e) and LEMON (Hu et al., 2022b).

In Encoder-Decoder architectures (Figure 2.3c), the encoder generates a cross-modal representation that is fed into a decoder that produces the final output. This architecture was originally adopted by LSTM-based models for text generation tasks and it has been popularized by Tranformers (Vaswani et al., 2017a) in Machine Translation and later in Natural Language Generation (NLG) by models like BART (Lewis et al., 2019) and T5 (Raffel et al., 2020). In VL the use of Encoder-Decoder architectures can enable the unification of VL understanding and generative VL tasks; examples of VL Encoder-Decoder models are OFA (Wang et al., 2022b), mPLUG (Li et al., 2022a) and MDETR (Kamath et al., 2021). A schematic representation of the architectures described is shown in Figure 2.3

<div align="center">(a) Encoder-only      (b) Decoder-only      (c) Encoder-Decoder</div>

Figure 2.3: Schematics of Transformers-based architectures. 2.3a Encoder-only with fully visible attention masking outputting a multimodal representations ($h^k$) the input consisting in the concatenation of visual tokens ($v^k$) and textual tokens ($t^k$). 2.3b Decoder-only using causal attention masking to autoregressively generate text starting from visual and optionally a textual tokens, e.g. a prompt. 2.3c Encoder-Decoder architecture combining both mechanisms to generate a text. Note that $N_V$ and $N_L$ are the total number of visual and textual tokens respectively. The figure is elaborated from Raffel et al. (2020).

## 2.2.5 | Pre-training Objectives

We can group the pre-training objectives into *intra-modality* and *inter-modality* objectives. The former enforces unimodal understanding and thus it helps the model to build relationships among the inputs of the same modality, though in some variants this process takes advantage of the other modality. The latter enforces interaction between textual and visual representations with the goal of building multimodal relationships.

The main *intra-modal* objectives are:

■ **Masked Language Modeling (MLM)**, introduced by Devlin et al. (2019a) as a textual pre-training task and then adopted also in VL pre-training. MLM consists in replacing a portion of the input tokens (Devlin et al. (2019a) originally replaced 15% of the total number of tokens) with a special [MASK] token. The model is optimized to reconstruct the input sequence by predicting the masked tokens with the original tokens. Dual stream models use this objective to learn robust textual rep-

resentations (e.g. CLIP (Radford et al., 2021b)), however in single-stream models (e.g. SOHO (Huang et al., 2021), UNITER (Chen et al., 2020b) and VisualBERT(Li et al., 2020c)), this task also exploits multimodal interactions, since the model has access to visual tokens and therefore the masked token prediction is conditioned on the visual information.

- **Masked Image Modeling (MIM)**, conceptually similar to MLM, but applied to the visual modality; it optimizes the model to reconstruct visual patches or regions given the remaining visible ones. There are some variations of this task: Tan and Bansal (2019) and Chen et al. (2020b) randomly replace the visual feature vector with zeros and train the model to regress the original feature with a mean squared error loss. Li et al. (2020e) and Lu et al. (2019), additionally, predict the object labels to which the regions correspond, leveraging labels obtained by an off-the-shelf object detector. BEiT (Bao et al., 2021a) uses MIM in combination with Vector Quantization (VQ) namely, the model is trained to reconstruct discrete visual tokens. A similar strategy is adopted in DALL-E (Ramesh et al., 2021) to perform text-to-vision generation, as this variation of MIM seems to help learning stronger visual representations.

On the other hand, the main *inter-modality* objectives in VL are:

- **Image-Text Matching (ITM)**, where given an image-caption pair, the model has to predict whether they match; this task is usually framed as a binary classification problem. Following Devlin et al. (2019a), most single stream VL models use the special `[CLS]` token prepended to the input sequence as a global input representation and feed it to a classification layer to predict a binary label (i.e. match/no-match). The native support for this task makes it a useful tool to perform systematic zero-shot evaluations on these models. We will go into detail on such techniques in Section 2.3.4. A notable variation of ITM is implemented in VisualBERT (Li et al., 2020c), where the model needs to predict whether an always-matching image-caption pair matches a second caption concatenated to the input. This implementation follows the Next Sentence Prediction (NSP) task originally used in Devlin et al. (2019a). **Word Region Alignment (WRA)** is a fine-grained version of ITM, introduced in UNITER (Chen et al., 2020b). The model has to optimize the alignment between image regions and words via the Optimal Transport algorithm. The aim is to minimize the cost of transporting the embedding distribution from image regions to words in a sentence. Chen et al. (2020b) shows that this method produces better joint embeddings for downstream tasks.

19

- **Image-Text Contrastive Learning (ITC)**, where given a batch of $N$ image-text pairs, the model is trained at the same time, to maximize the cosine similarity of the vector representations of the matching pairs and minimize the similarity between the non-matching pairs. Jointly optimizing the cosine distance of the unimodal representations forces encoders to project the textual and visual embeddings in the same vector space, creating de facto a multimodal joint vector space. This technique was conceptually known in NLP (Mikolov et al., 2013b) and initially introduced in VL by Cao et al. (2017). It was later popularized by CLIP (Radford et al., 2021b) and ALIGN (Jia et al., 2021) to pre-train large-scale dual stream VL models, as at least at large scale, it produces extremely robust and meaningful unimodal representations, which are comparable in the same vector space. Later on, it became a standard in VL pre-training; examples are ALBEF (Li et al., 2021b), BLIP (Li et al., 2022c), UNIMO (Li et al., 2020d), VLMo (Bao et al., 2022), CoCa (Yu et al., 2022), and FIBER (Dou et al., 2022) to name a few.

## 2.2.6 | How do VL models enforce multimodal grounding?

As discussed so far, we can identify two main mechanisms in VL modeling enabling multimodal grounding: (i) the *attention-based fusion module* (discussed in Section 2.2.3) that acts at an architectural level and enables explicit interactions at multiple levels between the modalities and (ii) via *objective optimization* (discussed in Section 2.2.5) where the model is forced to leverage multimodal information to optimize its parameters. However, several models implement only a number of these mechanisms in different settings and data. Thus, identifying the most effective ones in every condition is hard. Hendricks et al. (2021) directly address this issue by comparing VL datasets, objectives, and architectures in the same training conditions and evaluating on downstream tasks. They show that visual unimodal objectives like MIM do not improve the models' performance on downstream tasks. They also show that contrastive objectives like ITC are extremely beneficial to dual-stream models with no cross-modal attention (e.g. CLIP (Radford et al., 2021b) and ALIGN (Jia et al., 2021)). However, if the dual stream model is equipped with cross-modal attention and optimized on ITM, it does not gain any advantage from the contrastive loss optimization (i.e. ITC). Moreover, they find that dual-stream models with co-attention (e.g. LXMert (Tan and Bansal, 2019) and VilBERT (Lu et al., 2019)) perform roughly on par with single-stream models with self-attention (like UNITER (Chen et al., 2020b) and Unicoder-VL (Li et al., 2020a)) as also found by Bugliarello et al. (2021a).

## 2.2.7 | From task-specific to large-scale models

VL research has witnessed a gradual evolution marked by the exponential growth of model and dataset sizes. Gan et al. (2022) provides a historical overview, identifying two main stages marking the progressive transition from task-specific to large-scale pre-trained models. The main difference between these stages is in the increasing amount of data dedicated to the pre-training stage and in the constant scale-up of model sizes. This trend is motivated by findings on scaling laws governing large LMs training. Empirical observations suggest that models performance benefits from more data and parameters according to a predictable relationship (Chung et al., 2022; Kaplan et al., 2020). Although these observations were made on unimodal models, they inspired VL research to follow this direction.

**Small-scale task-specific models**   are an important line of work based on pre-extracted visual features (e.g. FasterRCNNs (Ren et al., 2015b) and ResNets(He et al., 2016)) and pre-trained textual features (e.g. Word2Vec (Mikolov et al., 2013a) and GloVe (Pennington et al., 2014)). These models were mostly based on Long Short-Term Memory (LSTM) (e.g. BUTD (Anderson et al., 2018)) for generative tasks such as image captioning and Fully Connected Network (FCN) (e.g. Visual Question Answering (VQA), (Antol et al., 2015)) for classification tasks. The idea was to reuse unimodal semantic representations from pre-trained models and successively train a decoder or a classifier on specific domain data. Curated datasets are usually used to pre-train unimodal feature extractors. Such datasets hardly exceed hundreds of thousands of samples, whereas the number of parameters typically does not exceed a million.

**Medium and Large-scale pre-trained models**   draw inspiration in VL by the success of BERT (Devlin et al., 2019a) in NLP and the successive adoption of Transformer-based models (Vaswani et al., 2017b) in CV research. The shift from medium (e.g. VisualBERT (Li et al., 2020c), LXMert (Tan and Bansal, 2019), OSCAR (Li et al., 2020e)) to large-scale pre-training, is mainly due to a consistent improvement of the performance of Pre-trained Vision and Language Models (PVL) on downstream tasks (e.g. CLIP (Radford et al., 2021b), Flamingo (Alayrac et al., 2022), SimVLP (Wang et al., 2021)) and the possibility to adapt larger models to downstream task with none (i.e. zero-shot and in-context learning) or only a few samples (e.g. few-shot learning), obtaining equal or superior performance to regular fine-tuning. Here, both the pre-training data and models' number of parameters increase massively, ranging from millions to billions for both. However, this approach produces noisy and biased datasets leading to serious risks in

terms of harmful content and biases in the models (Birhane et al., 2021; Schuhmann et al., 2022).

In the next Section (2.3), we will discuss the role of the data in current VL modeling, providing an overview of current tasks, datasets, and evaluation techniques.

# 2.3 | Data, Tasks, and Evaluation

## 2.3.1 | Pre-trainig Datasets

Data is the raw material to craft VL models. Large-scale models are trained on top of two stages: *pre-training* on general tasks (i.e. *pre-training objectives* as described in Section 2.2.5), and *fine-tuning* on downstream tasks. In the recent past, *pre-training* has become central in the development of large-scale VL models as it was shown that better pre-training can completely eliminate (zero-shot and in-context learning) or extremely reduce (one- or few-shot learning; ) the amount of data needed to adapt the model on downstream tasks (Fei-Fei et al., 2006).

Large pre-trained VL models have triggered a paradigm shift in the dataset creation process, which started from small **crowd-sourced** curated datasets to move towards large **web-crawled** noisy datasets. Below we report the main VL datasets surveyed; for an overview see Table 2.1.

**Crowd-sourced datasets** are small-scale datasets, commonly used in academic settings (Gan et al., 2022). The size of these datasets is in the order of a few hundred to hundreds of thousands of image-caption pairs, collected by tasking human annotators to describe images by emphasizing mostly the objects and the entities depicted in the image. This is enforced by means of specific guidelines that prevent the annotators from introducing any kind of additional subjective information (Hodosh et al., 2013b). Despite that, biases stemming from subjective interpretations can still be found as shown by Van Miltenburg (2016).

Our survey expands Ferraro et al. (2015) on crowd-sourced datasets. A popular example is the Pascal Dataset (Farhadi et al., 2009), one of the first datasets aligning images with textual descriptions; it contains $1k$ images each associated with 5 captions. The Flickr8k (Hodosh et al., 2013b) and its extension, Flickr30k (Young et al., 2014b), containing $159k$ crowd-sourced captions, describe people involved in everyday activities and events. Microsoft Common Objects in Context (COCO) dataset (Chen et al., 2015b) is one of the most popular VL datasets describing common objects and people in naturally occurring contexts. It contains 91 basic object types labeled in $328k$ im-

| | Pre-training Dataset | Images | Text | Scale |
|---|---|---|---|---|
| Crowd-sourced | Pascal (Farhadi et al., 2009) | 1k | 5k | S |
| | Flickr8k/30k (Hodosh et al., 2013b; Young et al., 2014b) | 8k/32k | 40k/159k | S |
| | COCO (Chen et al., 2015b) | 328k | 1.64M | S |
| | Déjà Images (Chen et al., 2015a) | 180k | 4M | S |
| | VG (Krishna et al., 2017b) | 108k | 5.4M | S |
| | LN (Pont-Tuset et al., 2020) | 628k | 650k | S |
| Web-crawled | CC3/12 (Changpinyo et al., 2021; Sharma et al., 2018c) | 3.3M/12.4M | 3.3M/12.4M | M |
| | SBU Captions (Ordonez et al., 2011b) | 1M | 1M | M |
| | WIT (Srinivasan et al., 2021) | 11.5M | 37.6M | M |
| | WenLan (Huo et al., 2021) | 30M | 30M | M |
| | *ALT200M (Hu et al., 2022b) | 200M | 200M | L |
| | LAION 400M/5B (Schuhmann et al., 2021, 2022) | 400M/5B | 400M/5B | L |
| | COYO 700M (Byeon et al., 2022) | 747M | 747M | L |
| | *CLIP (Radford et al., 2021b) | 400M | 400M | L |
| | *ALIGN (Jia et al., 2021) | 1.8B | 1.8B | L |

Table 2.1: Statistics for the small- (S), medium- (M), large- (L) scale pre-training datasets reviewed. We report the size in terms of *thousand/million/billion* (respectively with k/M/B) of image and text instances. All the datasets are publicly available except for ALT200M and those used in CLIP and ALIGN (denoted with *). Table adapted from (Ferraro et al., 2015).

ages, each associated with 5 captions. The Déjà Images Dataset (Chen et al., 2015a) is composed of $180k$ user-generated captions coupled with $4M$ Flickr images, by using a combination of rules, such as high-frequency nouns and human judgments collected on Amazon Mechanical Turk (AMT).

Visual Genome (VG) (Krishna et al., 2017b) is a densely annotated dataset, partially based on existing datasets, containing $108k$ images and 50 region caption annotations per image for a total of $5.4M$ region captions; making it probably the largest curated VL dataset, in terms of textual annotations. The dataset is provided with many different annotations like scene graphs, objects, attributes and relationship labels, and question-answer pairs to perform a number of different VL tasks. Localized Narratives (LN) (Pont-Tuset et al., 2020) is another remarkably large annotated dataset. It comes with $628k$ images coupled with $650k$ captions. Differently from previous works, the LN provides highly descriptive textual captions of images temporally aligned with mouse tracks pointing to the area described by the annotators.

Image-text pairs in crowd-sourced VL datasets have a high quality as they are collected in a controlled environment, and usually manually or semi-automatically quality checked during collection, however, they have a limited size usually due to high costs in terms of annotations, collection time, and quality monitoring. This makes them hard to scale, especially in limited resource settings like academic ones.

**Web-crawled datasets** are automatically collected by scraping from the web, allowing to easily collect a large amount of data at low or no cost. Web-crawling VL data has caught on in academic settings to boost VL research as it allowed the creation of medium-scale datasets (sized in the range of a few million image-text pairs). Conceptual Captions (CC) (Sharma et al., 2018c), for instance, is collected by scraping alt-text attributes and the corresponding images. The automatic pipeline included automatic filtering and heuristics to remove possible harmful content (e.g. pornography and profanity) and image-text pairs where there is little overlap between textual and visual content based on object labels extracted with an off-the-shelf OD. With this method, they collected $3.3M$ image-text pairs in the wild, subsequently extended to $12.4M$ (Changpinyo et al., 2021). SBU Captions (Ordonez et al., 2011b), contains $1M$ images coupled with user-generated captions, collected by performing semantic querying on Flickr similarly to Chen et al. (2015a). The captions are collected by querying for specific terms such as objects and actions.Wikipedia-based Image Text (WIT) (Srinivasan et al., 2021) a collection of $11.5M$ images and $37.6M$ alt-texts collected from Wikipedia content pages 108 languages. WenLan (Huo et al., 2021) is a dataset of $30M$ web-crawled image-text pairs. The sophisticated content-cleaning pipeline takes into account topic words and topic distributions.

This approach has gained interest also in industrial settings, giving rise to large-scale datasets (from hundreds of millions to billions of image-text pairs) (Gan et al., 2022). Examples of large-scale datasets are: COYO-700M (Byeon et al., 2022), containing 747M image-text pairs, collected following Sharma et al. (2018c) and LAION-400M/5B (Schuhmann et al., 2021, 2022), which contains $400M$ image-text pairs, later extended to $5B$. Differently from the previous datasets, the cleaning pipeline does not rely on handcrafted filtering rules, in fact, they use CLIP (Radford et al., 2021b) to automatically filter out poor matching image-text pairs, by computing the cosine similarity of their vector representations. However, the lack of proper curation can lead to datataset with potentially harmful contents as found by Birhane et al. (2021), analysing LAION-400M.

Other large-scale datasets not released to the public are: the dataset used in CLIP (Radford et al., 2021b), consisting of $400M$ image-text pairs, built upon a set of $500k$ queries containing terms occurring at least 100 times in the English Wikipedia. The dataset used originally in ALIGN (Jia et al., 2021), and later in SimVLM (Wang et al., 2021) and CoCa (Yu et al., 2022) is composed of $1.8B$ image-text pairs collected similarly to CC, but with looser cleaning steps. The cleaning pipeline consists of setting image-size limits, fixing alt-text frequencies, and excluding rare words. Hu et al. (2022b) uses a similar data collection pipeline collecting $200M$ image-text pairs, called ALT200M, used to train LEMON (Hu et al., 2022b), a state-of-the-art (SOTA) VL model. A recent ap-

proach, further expand the pre-training dataset by perfornibg bootstrapping with syn-
thetically generated data. This method called CapFilt (Li et al., 2022c) and further ex-
tended in (Li et al., 2023b) uses a captioning generator to produce alternative captions
which are then filtered by a filtering model trained separately.

## 2.3.2 | What kind of linguistic information is present in current VL datasets?

A common aspect of crowed-sourced and web-crawled datasets is the object-centricity
of the captions, namely captions that focus only on objects visible in the image. In
crowd-sourced datasets (e.g. COCO (Chen et al., 2015b)), this is enforced by instruct-
ing the annotators to explicitly mention only the visible content, in order to prevent the
introduction of any additional information that might come from subjective interpreta-
tion. In web-crawled datasets, this happens implicitly by using filtering methods that
favor captions mentioning objects and entities present in the image, like in CC (Chang-
pinyo et al., 2021; Sharma et al., 2018c). However, language is not only composed of
concrete words, such as objects and entities, in fact, as pointed out in Section 1.2, hu-
mans use language to express complex and abstract ideas. This kind of linguistic in-
formation is hardly found or It is not systematically highlighted in current curated and
web-crawled datasets. However, in the latter, we could observe its presence to a small
extent as medium- and large-scale datasets are collected with loose filtering pipelines. In
this work, we explore this direction. In Chapter 4 we try to shed light on the capability
of current VL models to handle and acquire such linguistic information.

## 2.3.3 | Downstream Tasks and Metrics

As described in Section 2.3.1, it is common practice to *pre-train* VL models on general
tasks (i.e. *pre-training objectives* described in Section 2.2.5) and then *fine-tune* and evaluate
on downstream tasks. VL research encompasses a broad range of tasks and applications.
In this section, we provide definitions, datasets, and metrics for the main tasks currently
studied in VL research. This is not intended to be a comprehensive review of VL tasks;
however, it aims to provide a brief overview of the main VL problems studied in the
literature. We follow Li et al. (2022b) and we group VL tasks into four categories, that
we describe below. For a schematic overview see Table 2.3.3.

**Retrieval tasks** are image-text matching (ITM) tasks. This is one of the fundamental
tasks in multimodality, indeed it is also used as a pre-training task for VL models (as

| | Task | Input | Output | Datasets | Metrics |
|---|---|---|---|---|---|
| Retrieval | Image-Text/Text-Image Retrieval | image/text | text/image | COCO (Chen et al., 2015b) Flickr30k (Young et al., 2014b) Flickr8k (Hodosh et al., 2013b) | Recall@K |
| Grounding | Phrase Grounding | image+text | bounding boxes | VG (Krishna et al., 2017a) Flickr30k (Young et al., 2014b) Flickr30k Entities (Plummer et al., 2015) | Recall@K, Accuracy |
| | Reference Expression Comprehension | image+text | bounding boxes | RefCOCO (Kazemzadeh et al., 2014) Talk2Car (Deruyttere et al., 2019) Visual7W (Zhu et al., 2016). | Accuracy |
| Visual Understanding | Visual Question Answering | image+text | text | VQA (Antol et al., 2015) | VQA Accuracy |
| | Visual Reasoning | image+text+graph | text | VCR (Zellers et al., 2019) CLEVR (Johnson et al., 2017). NLVR (Suhr et al., 2017b). NLVR2 (Suhr et al., 2018). GQA (Hudson and Manning, 2019). | Accuracy |
| | Visual Entailment | image+text | label | SNLI-VE (Xie et al., 2019) | Accuracy |
| Generative | Image Captioning | image | text | COCO (Chen et al., 2015b) Flickr8k (Hodosh et al., 2013b) Flickr30k (Hodosh et al., 2013b) SBU Captions (Ordonez et al., 2011b) | BLEU (Papineni et al., 2002a) ROUGE (Lin, 2004a) METEOR (Banerjee and Lavie, 2005) CIDEr (Vedantam et al., 2015) SPICE (Anderson et al., 2016) |
| | Text-to-Image | text | image | COCO (Chen et al., 2015b) CUB (Wah et al., 2011) | FID (Heusel et al., 2017), Inception score (Mao et al., 2016) |

Table 2.2: Overview of the main VL downstream tasks, metrics, and datasets surveyed in this section. Table adapted from (Li et al., 2022b).

described in Section 2.2.5). Given a query in a certain modality (i.e. textual or visual), the goal is to match or retrieve the closest element in the other modality. In VL this can be performed in two directions, namely *image-to-text* and *text-to-image*. Ranking metrics are used for evaluation, such as Recall@K. The main benchmarks for VL retrieval are COCO (Chen et al., 2015b), Flickr8k (Hodosh et al., 2013b) and Flickr30k (Young et al., 2014b).

**Grounding tasks** are designed the assess the capabilities of VL model to find connections between the two modalities. There are two main very similar tasks, namely, *phrase grounding* and *referring expression comprehension*. They both aim at identifying a bounding box in the image in relation to a given text, however in *phrase grounding* the goal is to find the region in the image that is better aligns with the text, whereas, in the *referring expression comprehension* the text is an expression that points to a specific subject or entity in the image, and the goal is to identify the region of the image the text refers to. Both tasks can be evaluated on an accuracy-based metric. The model predicts a bounding box for the input text: if the Intersection over Union (IoU) (Plummer et al., 2015) is greater than a specific threshold it is considered as a true positive. In *Phrase Gounding*, however, Recall@K is also used as the task can be framed as a ranking problem, where multiple boxes can be proposed as candidates for matching the input text. Common benchmarks for the *phrase grounding* task are VG (Krishna et al., 2017a), Flickr30k (Young et al., 2014b) and Flickr30k Entities (Plummer et al., 2015), whereas for the *reference expression comprehension task* the main benchmarks are RefCOCO (Kazemzadeh et al., 2014), Talk2Car

(Deruyttere et al., 2019) and Visual7W (Zhu et al., 2016).

**Visual Understanding tasks** have the ultimate goal of providing a system that can "reason" about the visual content. These tasks are typically framed as classification tasks and therefore the evaluation metrics are accuracy-based. Here we describe examples of such tasks.

*Visual Question Answering task* (Antol et al., 2015) where given an image-question pair, the model has to predict the correct answer among the possible ones; the evaluation is based on the accuracy metric proposed by Antol et al. (2015), based on human judgments. The task is proposed in two versions, namely multi-choice and open-ended generation. The assessment of the generative version is performed by converting gold and the model answers into bag-of-words. The answer is chosen to select the text among the choices available that overlap the most with the models' output. Then an accuracy-based metric is applied.

*Visual Reasoning task* is framed as a visual question-answering task, similarly to VQA (Antol et al., 2015), however, the model has to perform compositional reasoning requiring advanced semantic representations and thus the model is also provided with scene graph annotation to ease this process. For instance, in the Visual Commonsense Reasoning (VCR) (Zellers et al., 2019) benchmark the model has to infer actions, goals, and mental states. In the Compositional Language and Elementary Visual Reasoning (CLEVR) and Natural Language Visual Reasoning (NLVR) (Johnson et al., 2017; Suhr et al., 2017b) the model has to answer questions regarding objects' attributes and their relationships, such as spatial and logical operations in synthetic scenes. Cornell Natural Language for Visual Reasoning for Real (NLVR2) and GQA (Hudson and Manning, 2019; Suhr et al., 2018) are two visual reasoning datasets that expand this idea using real-world image captions. NLVR2 introduces a new task consisting of judging whether a caption is correct with respect to two images. GQA consists of automatically generated visually grounded questions leveraging scene graph annotations from VG. By using scene graphs annotations the authors where able to generate a large variety of questions, requiring very different reasoning skills.

In the *Visual Entailment task*, given an image-text pair, the model has to tell whether the image semantically entails, is neutral, or contradicts the input text. It relies on the notion of *logical entailment* and it is formulated as a classification problem. The SNLI-VE (Xie et al., 2019), is the main dataset to benchmark VL models on this task. Vu et al. (2018) provide a variation of this task called Grounded Entailment. In this setup, the premise is composed of image and text rather than only text, resulting in a more

difficult task as the model has to correctly ground both the textual information in the premise and the hypothesis.

**Generative tasks**   differ from classification tasks inasmuch as the output is not a numerical prediction, but is a signal in a specific modality (i.e. text or images). In VL, generative tasks can be performed in two directions: image-to-text and text-to-image. Differently from text-to-image, which has recently witnessed a surge of interest in the community (Frolov et al., 2021; Li et al., 2019; Ramesh et al., 2021), image-to-text, commonly known as *image captioning* is a classical VL problem as well as of primary interest for this work.

*Text-to-image generation* consists in generating an image from an input text. The evaluation is based on CV metrics, such as FID (Heusel et al., 2017) and Inception score (Mao et al., 2016) which measure the similarity of the image with the training distribution. This task is commonly assessed on VL datasets such as COCO (Chen et al., 2015b) and Caltech-UCSD Birds (CUB) (Wah et al., 2011).

The *image captioning task* consists in generating a textual description for an image given as input. Evaluating textual outputs is a well-known problem in NLG (Gatt and Krahmer, 2018; Howcroft et al., 2020; Sai et al., 2022). In *image captioning* the output evaluation often relies on n-gram-based automatic metrics such as BLEU (Papineni et al., 2002a), ROUGE (Lin, 2004a), METEOR (Banerjee and Lavie, 2005) and CIDEr (Vedantam et al., 2015), or on additional scene graph annotations like SPICE (Anderson et al., 2016), However, this is often combined with a human evaluation to obtain more consistent results. The main datasets used to benchmark models in the *image captioning tasks* are: COCO (Chen et al., 2015b), Flickr8k and Flickr30k (Hodosh et al., 2013b; Young et al., 2014b), SBU Captions (Ordonez et al., 2011b).

## 2.3.4 | Assessing Multimodal Grounding

While performance on downstream tasks is impressive, it is yet not clear what information models capture in their multimodal representations (Salin et al., 2022). Task-based evaluation is insufficient to assess the grounding capabilities of VL models as they can often solve tasks by exploiting artifacts and spurious correlations in the data (Ribeiro et al., 2016b; Wu et al., 2021; Yang et al., 2021, 2023). For example, in Natural Language Inference (NLI), where the model has to determine whether a premise entails, is neutral, or contradicts a hypothesis, the presence of a negation, such as "never" correlates with a contradiction in many NLI benchmarks (Gururangan et al., 2018; McCoy et al., 2019). Interpretable machine learning is a multidisciplinary field encompassing efforts from

computer science, human-computer interaction, and social science, with the goal of designing user-oriented and human-friendly explanations for machine learning models (Rudin et al., 2022).

Interpretable machine learning can be defined as the set of techniques, designed to extract the relevant knowledge concerning relationships learned by machine learning models. Such knowledge is considered relevant when it provides insights for a particular audience into a chosen domain (Murdoch et al., 2019). As a consequence, they help enhance the trust and confidence in these models. In this work, we focus on the application of interpretability techniques to validate and debug grounding capabilities models in the VL domain.

Interpretability techniques for Deep Neural Network (DNN) can be grouped into two categories: **white box** methods which exploit the knowledge of the internal structure of the model to generate the explanation and **black-box methods**, also called model-agnostic, which operate only on the inputs and the outputs.

Among the many **white box** methods, in this work, we focus on the following:

**Attention-based analyses**   rely on the analysis of attentional mechanisms (Clark et al., 2019b) to identify relations between multimodal inputs. Aflalo et al. (2022) proposes a comprehensive tool to analyze multimodal interactions in attention-based VL models. Similar analyses are often performed to analyze different attentional patterns and correlations between visual and textual tokens (Chen et al., 2020b; Li et al., 2020c,e; Tan and Bansal, 2019). However, while on the one hand it provides interesting qualitative visualizations, on the other, it is not clear how to aggregate and interpret information from deeper attention heads, though some methods have been proposed to track the attention flow across all the Transformers layers (Abnar and Zuidema, 2020). Moreover, it can be applied only to attention-based methods.

**Probing tasks**   are alternative approaches that allow inspecting information embedded in the multimodal representations. A probe is a classification model trained to solve a simple task using the information stored in the models' representation for the specific aspect tested by the task. The rationale behind this method is that a good probe performance is an indication of a meaningful representation. This approach has gained success in NLP, to test Transformer-based models (Penha and Hauff, 2020; Wallat et al., 2023) and subsequently adopted in VL. An example is the inter-modality probing proposed by Salin et al. (2022), which is used to quantify the interaction between vision and language modalities by training probes at intermediate layers, in order to be able to track the impact on the modality fusion across the model. Similarly, Cao et al. (2020)

design an extensive set of probes to test intermediate layer representations. Rather than focusing on the final output of the model, they use the attention layers as a proxy to observe intermediate cross-modal interactions. They find visual-explainable patterns revealing relations between the text segments and image regions.

In this work, we focus on the following **black box** methods:

**Behavioural benchmarks** are often employed to assess specific visio-linguistic capabilities. Shekhar et al. (2017) show that VL models lack fine-grained linguistic understanding when the caption is minimally edited, by replacing one of the mentioned entities with another one (*foiling task*). Along the same lines Hendricks and Nematzadeh (2021b) design a benchmark to test for verb understanding showing that VL models fail to ground verbs as they focus mostly on nouns, as previously observed by Tanti et al. (2018). This is in part due to the capability of the visual backbone to capture object-level information, as recently shown by Zhang et al. (2021). Using a similar method Thrush et al. (2022) design a task to test VL models for compositionality, whereas Parcalabescu et al. (2020) test models' counting capabilities. These tasks show that VL models perform poorly on tasks relatively simple for humans, highlighting a severe lack of multimodal grounding and the inadequacy of task-based evaluations. An example of a case study we contributed to, developed in parallel with this thesis, is VALSE (Parcalabescu et al., 2022a). VALSE is a benchmark designed for testing general-purpose pre-trained VL models on specific linguistic phenomena requiring visio-linguistic capabilities. The benchmark provides tests for five different linguistic phenomena (i.e. *existence, plurality, counting, relations, actions, coreference*) framed as a foiling task. The foils are automatically generated accounting for distributional biases, and extensively validated (automatically and with human validations).

**Input Ablation** is a technique consisting of measuring variations in the model's output by perturbing the input (Fernandes et al., 2021; O'Connor and Andreas, 2021). In multimodal settings, it can be used to analyze the impact of a modality-specific input or a portion of it. For example, Li et al. (2020e) performs an ablation study on the OSCAR model, to analyze the impact of object tags, added to the input to improve the model's capability to bridge the visual and the textual modality. Frank et al. (2021b) use a cross-modal input ablation to measure the extent to which VL models rely on one rather than the other modality.

**SHAP**   (Lundberg and Lee, 2017) is a framework considered by many to be the gold standard for local explanations, thanks to its solid theoretical background. SHAP leverages the concept of Shapley values, first introduced by Shapley et al. (1953), used to measure the contribution of players in a cooperative game. This was later extended by (Lundberg and Lee, 2017) for the purpose of explaining a machine learning model. Using the models' input features as players of a cooperative game where the models' output is the outcome of the game, SHAP estimates feature attributions, i.e. the Shapley values, that quantify the contribution of each input feature to the final output of the model.

In the scope of this thesis, we develop an automatic method based on *input ablation* (described in depth in Section 4.3), that combined with *zero-shot evaluations* allows us to analyse VL models behaviors based on the ablation of entities present in the visual and textual modality. In Section 4.4, we leverage *probing tasks* and extend *attention-based* analyses to understand the impact of the model's representations when exposed to particular linguistic information. Finally, in Chapter 5 we propose a flexible hybrid framework based on SHAP, which benefits from properties typical of *black-box* methods. At the same time our method shares features with *white-box* approaches since, when possible, it takes advantage of certain internal components of the model.

# 3

# Grounding actions, scenes and rationales: The HL Dataset

The material in this Chapter is based on: Michele Cafagna, Kees van Deemter, Albert Gatt. *HL Dataset: Visually-grounded Description of Scenes, Actions and Rationales*. Proceedings of the 16th International Natural Language Generation Conference (INLG2023);

**Contributions**: Michele Cafagna: collecting, processing and analysing the data; implementing and running the experiments; writing and revising the paper. Albert Gatt and Kees van Deemter: supervising the research, writing, and revising the paper.

| Image | Axis | Caption |
|---|---|---|
|  | scene | the picture is shot in a ski resort |
| | action | they are just relaxing after a round of skiing |
| | rationale | they want to have a good time together |
| | object-centric (COCO) | a woman and a boy sitting in the snow outside of a cabin. |

Table 3.1: Example of High-Level captions. It is shown one of the three captions available for the three axes collected: *scene, action, rationale*, combined with the object-centric captions from COCO.

# 3.1 | Introduction

Conceptual grounding broadly refers to the idea that symbols (e.g. language) are grounded in perception (Barsalou et al., 2008). As discussed in Section 1.2.2, perceptually grounded communication is made possible by the fact that perceptual experiences are largely shared. However, individual experience can also license subjective inferences which inform not just what we express through language, but also what we choose to assume and leave unexpressed (Bisk et al., 2020).

Among the many modalities available in the perceptual spectrum, visual grounding has always been of primary interest as it provides the easiest and most readily accessible way to link linguistic expressions to physical objects. As pointed out in Section 2.3.2, a glance at many widely used datasets and models in image captioning reveals a bias towards "object-centric" descriptions, whereby models are trained on image-text pairs where the text consists of explicit mentions of objects visible in the scene. However, experience and perception also motivate other, non-object-centric ways of talking about the world, for example, when we talk about scenes, or when we describe actions or their underlying rationales. While such "high-level" descriptions are also perceptually grounded, they incorporate world knowledge and subjective experience.

For example, the object-centric description in Table 3.1 certainly describes the visual content, though it is based mainly on the recognition of objects, and their spatial layout in the scene. By contrast, the three high-level captions (*scene, action, rationale*, from the HL-Dataset described below), provide three different perspectives of the scene among the many possible ones, which are triggered by expectations and assumptions based on subjective experience and world knowledge.

In this Chapter, we tackle the issue of grounding high-level linguistic descriptions in the visual modality, proposing the High-Level (HL) Dataset: a resource for Vision and Language (VL) modeling that aligns existing object-centric captions with human-collected high-level descriptions of images along three different axes: *scenes, actions* and *rationales*. The high-level captions capture the human interpretation of the scene which are complementary to object-centric captions used in current VL datasets, e.g. in COCO (Lin et al., 2014c). We took a step further, and we collected *confidence scores* from independent annotators, which serve to shed light on the extent to which the high-level captions in the dataset correspond to widely-shared assumptions, or to idiosyncratic interpretations. Finally, we considered the task of generating captions that incorporate these different axes, yielding a more narrative-like description of images. Our contributions are:

- We presented and released the HL Dataset, a new VL resource, grounding high-level captions in images along three different axes and aligned with existing object-centric captions;

- We described the collection protocol and provided an in-depth analysis of the data;

- We presented baselines for the High-Level Captioning task and described further potential uses for our data.

## 3.2 | Related work

As pointed out in Section 2.3.2 current VL datatets tend to focus on grounding objects and entities. As a result, the range of linguistic phenomena enriching VL models' grounding capabilities is significantly narrower then what potentially exploitable.

In this Section, we expanded on this issue, by analysing the historical reasons of object-centricity in VL and elaborating on the need for a change of focus that is in line with long-run goals initially set for the VL research.

Hodosh et al. (2013c), in their influential work, argue that image captioning is mostly interested in "conceptual descriptions", which focus on what is actually in the image and differ from the so-called non-visual descriptions, which provide additional background information. This line of thought has been broadly followed in the field, resulting in datasets emphasizing object-centric content in VL tasks involving text generation, like image captioning (Agrawal et al., 2019; Lin et al., 2014c; Sharma et al., 2018b) and visual question answering (Antol et al., 2015; Zhu et al., 2016). However, especially in crowd-sourced datasets where these conditions are explicitly enforced, annotators keep

making inferences based on their personal experiences (see e.g., Van Miltenburg (2016) for evidence). This supports the fact that humans find it natural to describe images by bringing to bear their knowledge and experience, using also more complex expressions conveying situational information e.g. a woman holding a child's hand is usually referred to as the child's mother or grandmother depending on the perceived age of the woman.

For instance, in the instructions used to collect COCO (Lin et al., 2014c), the annotators are explicitly asked to mention entities visible in the image. This is beneficial to enhance cross-modal interactions: Zhang et al. (2021) show that improving the visual backbone on object recognition tasks improves the performance of visio-linguistic models in downstream tasks. Li et al. (2020e) show that using object labels to bridge the two modalities improves the grounding capabilities of VL models. In a previous work, Wang et al. (2018) shows that in highly object-centric datasets such as COCO (Lin et al., 2014a) and FLickr30k Young et al. (2014b), image-captioning can be performed by replacing the visual information with object labels.

As discussed in Section 2.3.2, object-centricity is also a feature of widely-used web-scraped datasets: in the Conceptual Captions dataset for instance, Sharma et al. (2018b) filtered out all captions that did not overlap with object labels automatically identified by a computer vision model in the corresponding image.

Some efforts have been made to understand how low-level concepts improve generalization capabilities and connect to high-level concepts. Object-centric captions help to improve the generalization over unseen objects (Hu et al., 2021b) and play a role in the model understanding of abstract concepts (Cafagna et al., 2022; Wang et al., 2022c). In our work, we are interested in the relationships between what Hodosh et al. (2013c) refer to as "conceptual" and "non-visual" descriptions, which we re-frame as a distinction between low-level (object-centric) and high-level descriptions in multimodal learning. We released a novel dataset to foster research in this direction.

Motivation for the present Chapter is also provided by recent research exploring the visual correlates of inferences, temporal and causal relationships (e.g., Park et al., 2020), which also have implications for generation. In visual storytelling, for instance, a model has to understand actions and interactions among the visually depicted entities (Hong et al., 2023; Hu et al., 2020; Huang et al., 2016; Lukin et al., 2018). Identifying actions is a prerequisite for predicting their motivations or rationales as well as explaining automatically generated descriptions of images (Hendricks et al., 2018). Actions and intention are paramount to performing commonsense and temporal reasoning on visual inputs. Along these lines, Park et al. (2020) create dynamic stories on top of static images, where the task is to predict previous and subsequent events such as actions and rationales with

respect to a given situation depicted in the image. Our work is similar in spirit, as we align high-level descriptions of *actions* and *rationales* with low-level descriptions of static images.

Some work has also been done to test multimodal model grounding capabilities from a more linguistic perspective. Parcalabescu et al. (2022c) built a benchmark to test models on a variety of linguistic phenomena, like spatial relations, counting, existence, etc. Pezzelle et al. (2020b) assess the integration of complementary information of VL models across modalities, while Thrush et al. (2022) test multimodal models on compositional reasoning. In this context, the HL Dataset proposed here, can offer another benchmark for VL models' understanding of high-level descriptions of images, though in this Chapter we will focus only on the generative aspect. Such descriptions are licensed by the entities depicted in the visual modality and the relationships between them but they do not mention them explicitly.

## 3.3 | Data

In this section, we describe the protocol used to collect annotations for *scenes, actions* and *rationales* and the subsequent collection of confidence scores through crowdsourcing. Differently from previous works, such as COCO, where human annotators are instructed to be objective and to mention only the objects clearly visible in the picture, we elicited high-level captions by encouraging the annotators to rely on their subjective interpretation of the image.

### 3.3.1 | Data collection

The task of collecting high-level descriptions is by nature hard to define and requires a clear and careful formulation, therefore we ran a preliminary pilot study with the double goal of collecting feedback and fine-tuning the task instructions.

**Pilot**    The pilot was run with six participants who were trained on the task, with high proficiency in English and a background in computer science and linguistics. The participants were selected from our network (4 females and 2 males), with age ranging between $41 - 50$ for 5 participants and 1 in the range $18 - 30$.

With the feedback received from the pilot we designed a beta version of the task and we ran a small batch of cases on the crowd-sourcing platform. We manually inspected the results and we further refined the instructions and the formulation of the task be-

Figure 3.1: Annotation form presented to workers during the high-level captions collection. The instructions (shown in Figure A.1), are always visible to the annotators.

fore finally proceeding with the annotation in bulk. We show in Figure 3.1, the final annotation form.

**Annotation Procedure**    The participants were shown an image containing at least one human subject and three questions regarding three aspects or axes: *scene*, *actions* and *rationales* i,e. *Where is the picture taken?*; *What is the subject doing?*; and *Why is the subject doing it?* We explicitly asked the participants to rely on their personal interpretation of the scene and add examples and suggestions in the instructions to further guide the annotators. Moreover, differently from other VQA datasets like (Antol et al., 2015) and (Zhu et al., 2016), where each question can refer to different entities in the image, we systematically asked the same three questions about the same subject for each image. See Appendix A.1 for the full instructions and Appendix A.2 for details regarding the annotations costs. The annotation collection was ultimately performed by a total of 1054 annotators.

Figure 3.2: The confidence scores annotation form. We show the instructions, the image, the question, and the corresponding answer.

**Images**   As mentioned in Section 3.1 the COCO dataset has a very explicit object-centric orientation, therefore it provides a good starting point to select images, such that we could couple object-centric and high-level captions in a resource-lean approach. Moreover, the alignment of object-centric and high-level captions permits an investigation of the relationship between them.

We randomly selected 14,997 images from the COCO 2014 train-val split. In order to answer questions related to *actions* and *rationales* we needed to ensure the presence of a (human) subject in the image. Therefore, we leveraged the entity annotation provided in COCO to select images containing at least one person.

The whole annotation was conducted on AMT. We split the workload into batches in order to ease the monitoring of the quality of the data collected. Each image was annotated by three different annotators, therefore we collected three annotations per axis.

## 3.3.2 | Confidence Scores

The high-level descriptions were collected by asking the participants to interpret the scene leveraging their personal experience. The element of subjectivity led us to expect some variation in the resulting descriptions, especially where annotators needed

to infer actions and rationales. In order to distinguish what can confidently be considered widely-shared, or "commonsense" descriptions, from more idiosyncratic interpretations, we conducted a separate study where we crowd-sourced *confidence scores* for each high-level caption. We asked independent participants to score the likelihood of a high-level description given the image and the corresponding question on a Likert scale from 1 to 5. For a detailed example of the form see Figure 3.2.

**Agreement-based worker selection**   The confidence scores were collected following the same protocol used to collect the high-level descriptions. Using the data from our pilot study, which was carried out with participants who had been thoroughly briefed on the task, we ran a preliminary qualification task where we employed an *automatic worker selection method* to hire qualified annotators from the crowd-sourcing platform.

Let's consider the participants of the pilot as gold-standard annotators (as they were trained on the task) and their annotations as reference annotations. The inter-annotator agreement computed on the reference annotations can be seen as the gold-standard inter-annotator agreement $\alpha_{gold}$ of the task.

We ran the qualification task using the same set of items used in the pilot, then for each worker $w$ we re-computed the inter-annotator agreement (Hayes and Krippendorff, 2007), combining the workers and the reference annotations, obtaining $\alpha_w$. We computed an agreement ratio

$$r = \frac{\alpha_w}{\alpha_{gold}} \qquad (3.1)$$

Then, we selected a worker $w$ if $r > t$, where $t$ is a threshold empirically set to 0.5. This is equivalent to choosing workers such that their contribution does not negatively affect $\alpha_{gold}$ by a factor greater than $t$. In other words, the workers were selected if they were relatively compliant with the gold annotators. With this procedure we hired a total of 50 workers.

## 3.4 | Dataset Analysis

In this section, we analyse all the captions collected in the High-Level Dataset with the goal of exploring connections between our captions and the object-level captions already present in COCO. We analysed the distribution of the captions across different axes, also comparing them with the object-centric COCO captions[1].

---

[1]The analysis is performed by using Spacy v.3 pipeline for English using the `en_core_web_md` model to analyse the part of speech of the texts.

| Data | # Tok | Avg Len | # Uniq | # Cap |
|---|---|---|---|---|
| actions | 271168 | 6.02 | 7326 | 44991 |
| scenes | 233232 | 5.18 | 4157 | 44991 |
| rationales | 306396 | 6.81 | 8301 | 44991 |
| HL (tot) | 810796 | 6.00 | 12296 | 134973 |
| COCO | 857218 | 11.42 | 13300 | 75019 |

Table 3.2:  HL dataset caption statistics compared the COCO captions (object-centric) for the shared set of images.  We report the number of tokens (# Tok), average length (Len), number of unique tokens (# Uniq), and number of captions (# Cap).

## 3.4.1 | High-Level descriptions

We collected 3 annotations per axis over a set of 14,997 images for a total of 134,973 captions. An example of high-level descriptions aligned with the original object-centric caption from COCO is shown in Table 3.1, whereas Table 3.2 reports a more detailed comparison of the statistics. We expected to observe shorter texts in the high-level captions as annotators were not giving highly descriptive details typical of object-centric captions. This is visible in Figure 3.3, which shows that the length of the high-level captions is roughly half of the object-centric COCO captions. Though shorter, they have a comparable number of unique tokens (i.e. types) over all the axes (as reported in Table 3.2); this suggests that the high-level captions are not repetitive and contain a fair amount of lexical variability.

Moreover, as already mentioned, the COCO captions are object-centric, namely, these captions are collected to objectively represent the visual content. Although this is convenient in recognition-oriented tasks, they lack the situational knowledge required to contextualize scenes. Such knowledge is an essential part of the cognitive processes underlying the grounding of language in vision. Indeed, as shown in Figure 3.4, the most frequent lemmas in the original COCO captions for the images used in the HL Dataset denote mostly objects visible in the picture. The high-level captions represent the same visual content with the addition of situational knowledge coming from the three axes, and this is also visible in different lexico-semantic choices in the texts. For example, Figure 3.5 shows the most frequent lemmas found in the *scene* axis. Because we align them to the same images, the dataset gives us a clean way to explore the relationship between objects and high-level axes.

**Disentangling the content across the axes**    Asking the same three questions, i.e. *where- what-* and *why*, about the same main subject of the image allows us to consistently com-

Figure 3.3: Caption length of the HL captions divided per axis (action, scene, rationale) in comparison to the object-centric COCO captions (object).

pare the content of our captions across three well-defined axes. We analysed the most frequent nouns in the *scene* axis in order to characterize the kind of scenes mentioned in the captions collected. This will allow us to identify and quantify the distribution of scene types collected. The top most frequent scenes include *street, room* and *road*. These are scene types that can encompass a very broad variety of objects. However, we can also identify scenes for which a narrower range of objects is likely to occur with some regularity, for example, those related to sports activities like *baseball, tennis, ski, ground* and *court*, or domestic environments like *house, kitchen* and *living* (referring to 'living rooms'). For a more complete view see Figure 3.5 where we report the top 20 most frequent scenes in the HL dataset.

Similarly, we can characterize also the *action* and the *rationale* axes. We identified the *action* distribution by analysing the verbs contained in the captions. In Figure 3.6 we observe that the most frequent actions are related to sports activities, similarly to what observed in the *scene* axis distribution. The most frequent verbs are *play, ski, surf, skateboard*, but we can also find generic actions like *hold, walk, sit* and *eat*.

In the *rationale* axis we analysed both nouns and verbs. In this axis we expected to observe more subjectivity and content variability, with more lemmas denoting intents, mental states and events, including psych verbs. Our hypothesis is that the annotators leverage their personal experience to infer these answers to a greater extent than they

Figure 3.4: The most frequent nouns in the COCO captions of the shared set of images with the HL dataset.

Figure 3.5: The most frequent lemmas of the captions in the *scene* axis of the HL dataset.

do for scene descriptions.

The majority of the rationales express intentions; in fact, *want* is by far the most frequent term in the lemmas distribution. As observed with the other two axes, terms related to sports activities are more frequent (*play, game, tennis, practice*), but also related to leisure (*enjoy, fun, vacation, love, family*) along with generic activities (*work, wait, try, eat*). For more details see Figure 3.7.

The systematic disentanglement of the content along three axes can serve as a filter to identify or analyse sub-samples of the data with specific characteristics. For instance, as observed so far, we can confidently say that sports-related activities are predominant in the dataset.

**Connecting high- and low-level descriptions**   One of the main goals of this resource is to enable the discovery of connections between high- and low-level captions, which are descriptions of the same images at different levels of abstraction. By construction, the alignment provided by the HL Dataset allows us to identify concrete objects in images which provide "support" for inferring high-level concepts such as scenes, actions and rationales.

43

Figure 3.6: The most frequent verb lemmas of the captions in the *action* axis of the HL dataset.

Figure 3.7: The most frequent noun and verb lemmas of the captions in the *rationale* axis of the HL dataset.

We dive deeper into our analysis and study the connection between high-level concepts related to scene, action and rationale, to low-level objects present in the aligned COCO captions. We ask: "What are the most informative objects for a high-level concept (e.g. *enjoy*) found in a specific axis (e.g *rationale*)?"

We leveraged the Pointwise Mutual Information (PMI) (Church and Hanks, 1990) to find the most informative objects linked to a high-level concept. This was helpful to discover connections between concepts across different levels of abstraction but also gave clues on the content distributions within the axes. We filtered out object mentions which had a frequency less than 100 in the low-level captions. This left 475 object-denoting lemmas. Then, we computed the PMI between content words in the high-level captions and all these lemmas. For example, Figure 3.8 shows the nouns in the object-centric captions which have the strongest PMI with the verb 'enjoy' in the rationale axis.

We can observe that high-level captions can express different nuances of the same abstract concept. To take another example, *love* (in Figure 3.9) can refer to the love between an animal and its owner, between two partners (e.g. *wedding*) or the love for sports (e.g. *skate, snowboard*). In the same way, as shown in Figure 3.8 a general concept like *enjoy* can be characterized by object-level concepts leaning toward a specific nuance

44

Figure 3.8: Most informative objects for the word *enjoy* in the *rationale* axis. Font size is proportional to PMI.



Figure 3.9: Most informative objects for the word *love* in the *rationale* axis. Font size is proportional to PMI.

of meaning, like sports activities (e.g. *kite, snowboarder, skier*) or places (e.g. *sandy shore, ocean, lake*).



Figure 3.10: Most informative objects for the word *restaurant* in the *scene* axis. Font size is proportional to PMI.



Figure 3.11: Most informative objects for the word *kitchen* in the *scene* axis. Font size is proportional to PMI.

The PMI analysis provides interesting insight into the connection between object-level and high-level captions on all the three axes available.

On the *scene* axis, for instance, the PMI gives some clues on the extent to which an

Figure 3.12: Most informative objects for the word *look* in the *action* axis. Font size is proportional to PMI.

| Axis | Top Lemmas | Top Objects (PMI) |
|------|-----------|-------------------|
| scene | street | intersection, decker, meter |
| | room | living, wii, nintendo |
| | road | traffic, decker, intersection |
| action | play | nintendo, wii, swing |
| | ride | rider, carriage, wave |
| | hold | controller, remote, rain |
| rationale | want | mirror, bathroom, sink |
| | enjoy | wave, kite, ocean |
| | fun | wii, nintendo, controller |

Table 3.3: Top most informative objects of the top most frequent lemmas in the three axes (*scene, action, rationale*) according to PMI.

object can be considered diagnostic for a scene. For instance, two semantically similar scenes like *restaurant* (see Figure 3.10) and *kitchen* (see Figure 3.11) share several diagnostic objects, as we would expect. However, we can identify important semantic nuances: the scene *restaurant* contains objects related to the food (i.e. *pizza, cheese, wine, sandwich*) whereas *kitchen* contains objects related to the preparation of food (i.e. *stove, oven, tray, refrigerator*). Another example is shown in Figure 3.12, where the most relevant objects for the action *look* encompass a wide variety of contexts, like looking at a screen or a device (e.g. *device, screen, cellphone*) or entertainment (e.g. *zoo, zebra, giraffe*). For more examples see Table 3.3, which shows the top most relevant objects for the top three lemmas in the *scene, action* and *rationale* axes[2].

These semantic differences, while quite easy for humans to interpret, are not usually present in object-centric VL datasets. They are made explicit and easy to identify in the HL dataset, where captions with different levels of abstraction are aligned with the same image.

## 3.4.2 | Confidence scores analysis

Our confidence scores are similar in spirit to the *self-confidence* scores collected in the VQA dataset (Antol et al., 2015). However, they differ insofar as our scores are not self-

---

[2]Note that the PMI estimation is based on frequencies. This means that if an object has a relatively low frequency but co-occurs always with a kind of scene, action or rationale, its PMI with the respective scene, action or rationale will be very high. This is the case for objects like *nintendo* and *controller*.

Figure 3.13: Axis-wise confidence score distribution of the high-level captions.

reported by the authors of the captions but collected from independent annotators. The inclusion of an external judgment plays an important role in determining the reliability of interpretation operated by the annotators in the caption collection and therefore, in shedding light on the extent to which an annotator's interpretation of a scene relies on "shared" or "commonsense" knowledge, or is entirely idiosyncratic.

We observe an average confidence score of 4.47 on a Likert scale from 1 to 5 (with a standard deviation of 0.78 and a median of 5) over all the axes. This suggests that, overall, according to independent judges, our high-level captions succeeded in capturing shared or 'commonsense' high-level interpretations of the scene.

Furthermore, the confidence scores provide an additional perspective under which our data can be characterized: by performing an axis-wise analysis of the confidence scores distribution (see Figure 3.13), we observe that the *scene* and *action* captions feature the highest overall confidence, while the *rationale* axis lags behind by a small margin. We expect such differences, since determining the rationale of an action depicted in a static image is challenging, in particular, because annotators can leverage significant visual cues, but have no access either to temporal information or the subject's stated intentions. Therefore, they need to resort to their own priors and expectations which can also lead to idiosyncratic interpretations which independent judges – as in our confidence score analysis – would find relatively unlikely.

One important use of confidence scores is to provide a measure of uncertainty of the data, which can be used, for instance, to identify hard samples; an example is shown

| Idx | Scene caption | Confidence |
|-----|---------------|------------|
| 1 | in the restaurant | 1 |
| 2 | in the entrance of the library | 1 |
| 3 | the picture is taken outside a library | 3 |

Figure 3.14: Example of a "hard" sample in the HL dataset where the scene captions have low confidence scores.

in Figure 3.14. The scene is hard to interpret even for humans and the scene captions display more variability and have low confidence scores.

### 3.4.3 | Quantifying Lexical and Semantic Diversity

In Section 3.4.2, we showed that in the presence of low confidence, there can be variation or disagreement among high-level captions given by different annotators for the same axis. In such cases, the captions focus on different aspects or refer to different interpretations. Although this phenomenon has been observed for captions with a low confidence score, it is conceivable that it might also happen with high-confidence captions, for example, two captions annotated by different annotators, while differing in the interpretation of an image, could nevertheless be considered highly likely. To quantify this phenomenon, in this section we further expand our analysis by studying the lexical and semantic diversity of our captions.

**Purity score**    We leveraged the BLEURT score (Sellam et al., 2020), a trainable metric used to evaluate semantic differences in NLG, to compute a score measuring the semantic diversity among the high-level captions associated with an image. To do so, we first computed such scores across each axis, and then we combined them to obtain a final score for the item. In this way, we could unpack the semantic diversity item-wise and axis-wise.

Let $C$ be the set of high-level captions of a given axis (e.g. scenes) for a given image. For simplicity, we do not report the index of the image and the axis in the following notation. Given a caption $c_i$, $\forall\, i = 0..|C|$, we computed its purity score $s_i$ as follows:

$$p_i = BLEURT(c_i, ref) \tag{3.2}$$

Here, $ref$ is the set of reference captions defined as:

$$ref := \{c_j \mid c_j \in C \text{ and } j \neq i\} \tag{3.3}$$

In practice, since for an image in a given axis in the HL dataset we have three captions i.e. $|C| = 3$, the purity score of a caption is the BLEURT score computed using the other 2 captions as references. $p_i$ gives a measure of the semantic diversity of the caption with respect to the other captions along the same axis.

By averaging the purity scores of all the captions across a single axis and across all the axes we obtained respectively a *purity score* measuring the semantic consistency both axis-wise and item-wise.

**Diversity score**    Along the same lines, we propose the *diversity score*, to measure the lexical diversity of the captions. The *diversity score* follows the same logic implemented to compute the *purity score* introduced in the previous paragraph, but the BLEURT score in Eq. 3.2 is replaced by the BLEU score (Papineni et al., 2002b) and then normalized between 0 (similar) and 1 (very different). This is achieved by performing a min-max normalization to the BLEU scores and then inverting the normalized score, namely computing $d_i^{final} = 1 - d_i^{norm}$. Here $d_i^{norm}$ is the min-max normalized BLEU score computed.

Our score is similar in spirit to self-BLEU (Zhu et al., 2018) as it measures the similarity of the captions within their own distribution. However, its computation concerns only axis-wise and item-wise captions.

**Results and discussion**    As shown in Figure 3.15 the purity scores obtained are mostly negative, this is due to lexical variations, which the BLEURT score is known to be sen-

sitive to (Sellam et al., 2020). However, BLEURT is not defined in any specific interval thus, it is usually hard to interpret (Sellam et al., 2020) if not considered in relative terms.



Figure 3.15: Axis-wise purity score distribution.



Figure 3.16: Axis-wise diversity score distribution. The scores have been normalized between 0 and 1.

Based on that, we use it to compare the semantic purity across items and axes within our dataset. As shown in Figure 3.15, *action* and *scene* share similar purity score distributions whereas the *rationale* is more skewed to the left than the scenes and *actions*. This shows that the rationales feature a higher semantic diversity (lower overall BLEURT) than the other axes.

The *rationale* axis is also the one featuring the highest lexical diversity, whereas the *scene* and the *action* have similar distributions. This is shown in Figure 3.16 where the *rationale* density estimate (in green) has a higher peak skewed on the right-hand side than *scene* and *action* density estimate (respectively in orange and blue).

We have similar observations for both *purity* and the *diversity* scores and this confirms what was observed in the confidence score analysis in Section 3.4.2, namely that the task of determining the rationale of an action from a static image produces more variation and divergent interpretations leading to higher semantic and lexical diversity. This is partly confirmed by the general observation that, *purity* and *confidence* scores positively correlate with each other, whereas *diversity* has a slight negative correlations with the two scores (See Figure 3.17).

For more details on the item-based analysis see Appendix A.3.

Figure 3.17: Pearson correlation between confidence, diversity and purity scores.

| Original | Review | Type of Error |
|---|---|---|
| it wants to take rest | it wants to take a rest | prepositions/articles |
| he eat | he eats | verb conjugation |
| is travelling to a particular place | he is travelling to a particular place | pronoun omission |

Table 3.4: Examples of the most common errors found by the annotators in the HL dataset's captions. We highlight in red the wrong part of the original captions collected and in blue the part corrected or added by the annotators.

## 3.4.4 | Quantifying grammatical errors

Despite constant monitoring of the annotation quality during data collection, we also performed a post-hoc assessment of the grammatical quality in order to evaluate the reliability and validity of the annotated data and identify potential sources of error. We asked two Master students in linguistics to correct grammatical errors in a sample of 9900 captions, 900 of which are shared between the two experts. They were shown the image-caption pairs and they were asked to edit the caption whenever they identify a grammatical error. The most common errors reported by the annotators are:

■ Misuse or lack of prepositions/articles

■ Wrong verb conjugation;

■ Pronoun omissions.

Examples of each error are shown in Table 3.4.

In order to quantify the extent to which the corrected captions differ from the original ones, we computed the Levenshtein distance (Levenshtein, 1966) between them.

We observe that 22.5% of the sample has been edited and only 5% with a Levenshtein distance greater than 10. This suggests a reasonable level of grammatical quality overall, with no substantial grammatical issues. This can also be observed from the Levenshtein distance distribution reported in Figure 3.18. Moreover, the human evaluation is quite reliable as we observe a moderate inter-annotator agreement ($\alpha = 0.507$, (Krippendorff, 2018)) computed over the shared sample.



Figure 3.18: Distribution of the Levenshtein distance computed between the original and the corrected high-level captions in a sample of 9900 captions.

## 3.5 | Generating high-level captions

In this section, we show how the dataset can be used to fine-tune models to generate high-level, aspect-specific descriptions, e.g. image-to-scene or image-to-action. Below, in Section 3.6, we also describe a data augmentation and generation experiment, to merge the three axes into more "narrative-like" descriptions of images.

We split the HL dataset reserving 90% for training and 10% for testing, accounting respectively for 13498 images and 121482 captions for training and 1499 images and 13491 captions for testing. In our experiments we reserved 10% of the train set for val-

| Model | Axis | CIDEr | SacreBLEU | Rouge-L |
|---|---|---|---|---|
| GIT | action | 110.63 | 15.21 | 30.43 |
|  | rationale | 42.58 | 5.90 | 18.57 |
|  | scene | 103.00 | 24.67 | 33.92 |
| BLIP | action | 123.07 | 17.16 | 32.16 |
|  | rationale | 46.11 | 6.21 | 19.74 |
|  | scene | 116.70 | 26.46 | 35.30 |
| ClipCap | action | **176.54** | **27.37** | **39.15** |
|  | rationale | **78.04** | **11.71** | **25.76** |
|  | scene | **145.93** | **36.73** | **42.83** |

Table 3.5: Automatic metrics for baselines (GIT, BLIP, and ClipCap) fine-tuned along the three axes (*scene, action*, and *rationales*) of the HL dataset. The results are the average of 5 evaluation runs, by keeping the same decoding strategy and parameters for all the models.

idation. We provided baselines for this task by fine-tuning three models, namely GIT (Wang et al., 2022a), BLIP (Li et al., 2022c), and ClipCap (Mokady et al., 2021) on each separate axis. In the choice of the baselines models we try to find a good trade-off between performance on one hand and size and efficiency on the other. Below, we give an overview of the models. All the baselines were fine-tuned for a maximum of 10 epochs using a learning rate of $5e-5$, Adam optimizer, and half-precision (`fp16`).

**GIT**   (Wang et al., 2022a) follows a standard and simple approach by employing a Transformer-based encoder-decoder architecture, not relying on any external OD or Optical Character Recognition (OCR) model. The model uses a pre-trained image encoder to generate visual tokens which are concatenated to textual tokens. The whole sequence is fed into a textual decoder which is trained from scratch using a causal attention mask to generate text. GIT is trained on a mix of open source and web-scraped data for a total of $0.8B$ image-text pairs using only the MLM objective. The model accounts for a total of $700M$ parameters.

**BLIP**   (Li et al., 2022c) uses a quite sophisticated training schema. An image and a text encoder are jointly optimized using an ITC objective. Moreover, an image-grounded text encoder and decoder are trained injecting visual information from the unimodal image encoder using cross-attention. While the first is optimized on ITM, the second is optimized on a LM objective. The pre-training dataset is bootstrapped by using a method called CapFilt, consisting in using and image captioner to generate synthetic captions and a filter model which selects the best candidates. This allows to improve

the quality of the pre-training dataset which accounts to $129M$ images and $130M$ texts. In spite of its complexity, this Vision and Language Pre-training (VLP) allows to train four re-usable models (i.e. image encoder, text encoder, image-grounded text encoder and image-grounded text-decoder) applicable to a wide variety of multimodal and uni-modal tasks. Moreover, the heavy use of parameter sharing across the models, during training, allows to keep down the total number of parameters, i.e. $252M$.

**ClipCap** (Mokady et al., 2021) successfully re-uses a pre-trained LM for captioning by exploiting prefix-based training. The method consist in training a simple mapping network to project the visual representation from a pre-trained image encoder, i.e. CLIP (Radford et al., 2021b), to a prefix to be fed into a pre-trained Large Language Model (LLM). The LLM can be frozen or optimized along with the mapping network. This method is simple, efficient and produces decent results when compared to more complex and heavy-compute SOTA solutions.

**Discussion:** Table 3.5 displays automatic evaluation results for the three models, on each axis. The first observation is that ClipCap outperforms by far the other models in each separate axis. Differently from the other models, which are natively multimodal, ClipCap leverages a LLM to generate captions, conditioning the text generation on a prefix representing the visual information.

A second observation, consistent with the analysis presented in earlier sections, is that over all the metrics, models fine-tuned to generate rationale-based descriptions receive lower scores than the other axes. We hypothesise that this is due in part to its inherent difficulty, as reflected in lower confidence scores. As shown in Figure 3.19, metrics obtained by models trained on the *rationale* axis, have smaller correlations with *diversity* and *purity* scores than models trained on the other axes, i.e. *scene* and *action*, while correlating similarly in terms of *confidence* scores. This partially confirms our hypothesis on the difficulty of the task of inferring the rationale of a situation happening in a scene. Indeed, lexical (*diversity* score) and content variety (*purity* scores) have a less systematic covariance relationship on the rationale generation (lower correlation) then on the other axes, despite having lower performance on the evaluation metrics. Future work could leverage these scores as additional signal in fine-tuning models on captions that require more inference, compared to more descriptive ones.

54

Figure 3.19: Average Pearson correlation between the HL dataset scores (confidence, diversity and purity) with the evaluation metric results (SacreBLEU, ROUGE, CIDEr) obtained by averaging fine-tuned models results axis-wise across the single items.

| Manually Annotated | T5-generated |
|---|---|
| sitting as a group with colleagues inside an office to take a group photograph | The group can be sitting within an office together with colleagues to take a group photograph. |
| He is surfing at the beach because he is trying to have fun. | he is surfing on the beach he is trying to have fun. |
| she is lying on the snow on the skating field because she has lost his balance and fell | He is lying on the snow on the Skating Field and has lost balance & fell. |

Table 3.6: Examples of manually annotated (narrative-like captions) obtained combined the three axes, i.e. *scene*, *action* and *rationale* and T5-generated captions obtained with the human-in-the-loop fine-tuning.

## 3.6 | Narrative-like generation

We now describe how we extended the dataset to combine the three axes to compose a short "narrative", which describes the scene, action and rationale in tandem. We called this new dataset HL Narratives. To do this, we leveraged the individual axes and synthesise this part of the data using a pre-trained language model. Since scenes, actions, and rationales were elicited individually in a visually grounded and controlled setting, a synthesised version of the three individual captions should also match the image to the same extent (modulo the variations in confidence that we observe).

## 3.6.1 | Data generation process

We framed the synthesis of narrative-like captions as a paraphrasing task. We used two alternative approaches: *human-in-the-loop fine-tuning* and *few-shot prompting*.

We tested the following generation approaches:

**Approach 1: Human-in-the-loop fine-tuning**  We followed a human-in-the-loop approach consisting of three stages: (i) we manually annotated a small sample of gold data; (ii) we fine-tuned a LLM on our gold data; (iii) we used the fine-tuned model to generate a sample of data, which is manually corrected and then (iv) added to the gold annotations before fine-tuning again. This procedure allowed us to use only a few iterations to annotate quickly a considerable amount of data. Since the model improves the quality of the generated data, it made the manual correction progressively easier to perform.

We used a version of T5 (Raffel et al., 2020) already fine-tuned on paraphrase generation[3] as our LLM data generator. We initialised the process with manually paraphrased annotations for 50 images ($3 \times 50 = 150$), fine-tuned the model for 2 epochs, and generated 150 captions for another 50 images, which were manually corrected and added to the original 150. The model was then fine-tuned for further two epochs. In each iteration, we reserved 10% as validation data. After two epochs, we observed that the validation loss did not improve further. Finally, in the last iteration, we used all gold data to fine-tune the model and generate synthetic high-level captions for the whole HL dataset, obtaining 14,997 synthetic captions for training and 1499 for testing. See Table 3.6 for some examples of manually annotated and T5-generated narrative-like captions.

**Approach 2: Few-shot prompting**  We built a data generation pipeline by leveraging the in-context learning capabilities featured by the most recent LLM (Brown et al., 2020a; Maeng et al., 2017; Touvron et al., 2023). This data generation approach has the advantage of not requiring any model fine-tuning.

We designed a prompt for our task and we used it to generate data from the recently developed LLaMA model (Touvron et al., 2023). The prompt consisted of the task description, followed by an example and the inputs of the task written in natural language. The full prompt is shown in Figure 3.20. The resulting output was then post-processed to extract the generated high-level caption.

---

[3]Details about the T5 fine-tuned on paraphrase generation are available at `https://huggingface.co/Vamsi/T5_Paraphrase_Paws`.

> Given three sentences merge them into one sentence, and make sure that the sentence is grammatically correct. Here is an example:'in a beach',' holding an umbrella',' so they won't get a sunburn' <holding an umbrella in the beach so that they won't get a sunburn.>\n The three sentences are: <**'scene'**,**'action'**,**'rationale'** >

Figure 3.20: Prompt used for the data generation. The parts in bold are replaced with the corresponding high-level descriptions for the given sample.

| Model | SacreBLEU | ROUGE-L | Cider |
|---|---|---|---|
| GIT (PRE) | 1.23 | 11.91 | 18.88 |
| GIT (T5) | **11.07** | **31.37** | **74.79** |
| GIT (LLaMA) | 10.96 | 24.71 | 65.05 |

Table 3.7: Automatic metrics computed over the gold annotated high-level captions; the scores are the average results of 5 runs using the same decoding parameters for all models. We compare the pre-trained model (PRE) with the model finetuned on T5-generated (T5) and LLaMA-generated (LLaMA) data.

**Synthetic data selection** We evaluated the quality of the two data generation approaches by comparing the output of the same baseline model, when it was fine-tuned on synthetic data which was generated with the two different methods. The goal was to determine the best synthetic dataset to fine-tune the baselines on. We used GIT-base as baseline image captioning model and fine-tune on the LLaMA- and T5-generated synthetic data. We evaluated the two versions of the model on a combination of qualitative models output inspections and automatic metrics (SacreBLEU (Post, 2018), ROUGE-L (Lin, 2004b) and CIDEr (Vedantam et al., 2015)) computed against the gold data.

In Table 3.7, we show the results of the evaluation based on the automatic metrics. First, we observe that the performance of the pre-trained model (PRE) is extremely poor in the high-level caption generation task, highlighting the substantial difference between captions of this kind with traditional object-centric captioning the pre-trained model is trained on. Some examples are shown in Figure 3.21.

Second, focusing on the fine-tuned models, we observe that GIT fine-tuned on T5-generated data performs better than the LLaMa-based counterpart on the automatic metrics. We argue that the model trained on T5-generated synthetic data benefits from the exposure of the model to the gold data distribution. However, we point out that the few-shot data generation pipeline remains a valid alternative as it achieves comparable performance without requiring any further fine-tuning.

**GIT (PRE)**: a group of people on the beach
**GIT (FT-T5)**: people enjoying sunbathing, the picture was taken on the beach and are going to have fun and entertainment

**GIT (PRE)**: two girls looking at their cell phones
**GIT (FT-T5)**: they are reading a text message outside on the street, waiting for their friend.

Figure 3.21: Comparison between the object-centric captions generated by GIT pre-trained (PRE) and the high-level caption generated by the fine-tuned (FT-T5) model, on T5-generated narrative-like captions. The generated high-level caption embeds high-level information regarding action, rationale, and scene, depicted in the visual content.

| Model | SacreBLEU | ROUGE-L | Cider |
|---|---|---|---|
| GIT (PRE) | 1.23 | 11.91 | 18.88 |
| BLIP (PRE) | 3.47 | 15.21 | 24.15 |
| ClipClap (PRE) | 8.72 | 19.45 | 40.47 |
| GIT (FT) | 11.11 | **27.61** | 75.78 |
| BLIP (FT) | **11.70** | 26.17 | **79.39** |
| ClipCap (FT) | 8.15 | 24.53 | 63.91 |

Table 3.8: Results of the narrative generation task, averaged over 5 runs using the same decoding parameters for all models. PRE: pretrained models; FT: finetuned on the synthetic data.

## 3.6.2 | Results

We built three baselines by fine-tuning the same three large pre-trained models used in Section 3.5: GIT, BLIP, and ClipCap on our T5-generated synthetic narrative-like captions. We fine-tuned for 3 epochs with batch size 8, learning rate $5e^{-5}$, and Adam optimizer with weight decay (Loshchilov and Hutter, 2017). We tested on our gold human-annotated data.

**Automatic metrics**    As shown in Table 3.8, where we report results for automatic metrics, overall the models achieve worse results than in the aspect-specific caption generation task (reported in Table 3.5). This further highlights the difficulty of generating narrative-like captions of this kind for models trained on object-centric captions.

**GIT**: he is riding a bike in the woods and is going to work.
**BLIP**: he is riding a bike in a park he is going to work.
**ClipCap**: He is riding a motorcycle on a road, he is riding a motorcycle because he wants to enjoy the ride.

**GIT**: the dog is jumping in the air, the picture is taken in a park, he is jumping
**BLIP**: the dog is jumping in the air. the picture is taken in a ground and he is doing
**ClipCap**: He is riding a skateboard in the snow, he wants to get to the top of the mountain.

Figure 3.22: Randomly picked examples of narrative-like captions generated by our baselines fine-tuned on the synthetic data.

Notably, the best-performing model in the aspect-specific caption generation task, namely ClipCap, is the worst performing model in the narrative-like caption generation, though by a small margin (Table 3.8). This suggests that although a conditioned LLM can greatly adapt to generate high-level descriptions of specific aspects of the scene, it struggles in generating comprehensive high-level descriptions involving multiple high-level aspects of the scene. Ultimately, this suggests that the multimodal representations learned by multimodal models are more robust and effective in generating natural captions than conditioned unimodal models such as ClipCap.

However, the exposure to a small amount of synthetic high-level captions is sufficient to drive the models' generated text toward more narrative-like outputs.

**Qualitative assessment**   We manually inspected a batch of 100 randomly picked narrative-like captions generated by our baselines. We performed a systematic error analysis, by classifying hallucinations or errors found in the captions into five categories:

1. subject hallucination, namely hallucinations regarding the main subject of the scene, such as misuse of pronouns;

**Error Type**: Subject
**Caption**: he is cooking in a kitchen he is hungry

**Error Type**: Action
**Caption**: they are discussing their company in a newly married building, they need to understand the project.

**Error Type**: Scene
**Caption**: they are taking a photo in an airport they are on a trip

**Error Type**: Repetitions
**Caption**: the picture is taken in a zoo, the zoo is a zoo and the zoo is a zoo...

**Error Type**: Other
**Caption**: the sun is visible because it is a beautiful day in the city

Figure 3.23: Examples of hallucinations found in the baselines' output for each error type.

2. action hallucinations, meaning that the action described does not reflect the actual action depicted in the image;

3. scene hallucination, concerning the wrong attribution of the location depicted in the image;

4. word or sentence repetitions;

5. other kind of hallucinations or errors not falling in the above categories.

In the assessment of the presence of the hallucination, we took into account both the image and the gold annotations. As already observed in Section 3.4 some axis-wise

Figure 3.24: Type of hallucinations and errors observed in a sample of 100 manually inspected narrative-like descriptions generated by our baselines (BLIP, GIT, ClipCap). We also report report the number of captions for which no error is observed (marked as 'none'). Note that more then one error may be observed in a caption therefore the total count may not add up to 100.

descriptions differ for the same image though being equally plausible. Therefore, a hallucination was found when the caption contained information that was not mentioned or anyway present in the gold annotations. Additionally, we excluded the *rationale* information in the narrative-like description from consideration as a possible source of hallucinations because, as already argued in Section 3.4, rationales are very subjective and most of that information cannot be unambiguously verified. In Figure 3.23 we show some examples of errors found.

As shown in Figure 3.24, ClipCap's generations feature a significant number of hallucinations if compared with the other two models. This observation is consistent with the automatic metric results shown in Table 3.8, where ClipCap the worst performing baseline. Interestingly, the majority of the hallucinations observed in this model involve mostly repetitions. If we consider only hallucinations concerning *scene* and *action* axes, ClipCap produces the lowest number of hallucinations. Moreover, it is the only model which does not produce any hallucination regarding the subject of the scene.

The results for this model are in contrast to what is observed in the axis-wise caption generation performed in Section 3.5, where ClipCap was the best performing model.

This suggests that differently from native multimodal models, such as GIT and BLIP, ClipCap struggles to adapt to the narrative-like captioning style. However, it is unclear why this is observed only with this kind of description. We speculate that this may be related to the mechanism employed in ClipCap to adapt a unimodal model to process the visual modality. Learning a prefix embedding to inform the language model on the visual information may produce an information bottleneck, leading to an unexpected output degradation. However, we also point out that the use of GPT-2 (Radford et al., 2019) as textual decoder in ClipCap may have some limitation in its own right on longer, narrative-like texts.

GIT and BLIP generate overall better quality narrative-like descriptions then Clip-Cap, with BLIP being the best performing model (Figure 3.24). They both produce a higher number of hallucinations than ClipCap, most of which concern the subject, the scene and the action, with the latter being the most common hallucination found in these models.

In Figure 3.24 we can observe that the number of correct captions (marked as 'none') in GIT is similar to ClipCap. This may seem in contrast with the automatic metric results reported in Table 3.8. However, the high rate of repetitions and unclassified hallucinations present in ClipCap's outputs, has a higher impact on the general quality of its generation outputs, resulting in overall better results in the automatic metrics for GIT rather than ClipCap.

In Figure 3.22 we show two randomly picked examples. They fairly describe the three axes (*scene, action, rationale*) in a single caption, though with some limitations in terms of fluency, as the axis-wise information is often not properly connected. In fact the axis-wise information in the narrative-like descriptions is oftentimes tied together by punctuation and conjuctions. This causes the overall degradation of the perceived quality of the narrative. However, we believe that this is in part due to the limited quality of synthetic data that does not match the quality of the human annotated captions collected for the single axes. See Appendix A.4 for more examples from all models.

Further progress can be done in this direction, for example by incorporating confidence scores during fine-tuning or extending the size of the manually annotated narrative-like descriptions.

## 3.7 | Further uses of the HL Dataset

We envision a wide set of further use cases and tasks enabled by the HL Dataset.

**VL generative tasks**   Our captions support image captioning generation tasks which encompass a broader range of visually grounded linguistic descriptions than the highly object-centric, "conceptual" descriptions that dominate the captioning literature Hodosh et al. (2013c). Moreover, the decomposition along three axes can be exploited to compose narratives of the image, as in image paragraph generation (Wang et al., 2019) and visual storytelling (Hu et al., 2020; Huang et al., 2016). They can be used in combination with the question each axis corresponds to, in order to generate micro-dialog scenarios.

We would also argue that the high-level captions are also more natural and human-like, since they were collected without enforcing any restriction on the content to be described. Given that the images are also aligned with object-centric captions, it is possible to envisage a scenario in which a model is trained to generate high-level captions, which are "explained" or justified with reference to low-level, object-centric properties (see Hendricks et al., 2016, 2018, for some work in this direction). In this way, the dataset can be leveraged to provide captions and explanations. Furthermore, the confidence scores serve for the identification of hard samples in the data, both for evaluation purposes and to provide additional training signals, as recently shown by Ouyang et al. (2022).

**Multimodal Grounding**   HL Dataset is also a useful resource to benchmark the grounding capabilities of large pre-trained VL models. Along these lines, Cafagna et al. (2021) study the capability of VL models to understand scene descriptions in zero-shot settings, finding that only large-scale pre-trained VL models have enough generalization capabilities to handle unseen high-level scene descriptions. Cafagna et al. (2022) analyse the impact of exposure to high-level scene descriptions on multimodal representations in models pre-trained on object-centric captions. They show that exposure to high-level concepts mainly affects the model's attentional resource allocation over the visual input, even though the low-level concepts learned during pre-training provide enough signal to support and easily adapt to scene descriptions during fine-tuning. This is also supported by Wang et al. (2022c) who find that low-level concepts are needed to learn higher-level concepts, though this does not hold in the other direction.

## 3.8 | Summary

In this Chapter, we introduced the High-Level (HL) Dataset. We extended 14,997 images from the popular COCO dataset with 134,973 human-annotated high-level descriptions

systematically collected over three axes: *scene*, *action*, and *rationale*. We aligned high-level captions with object-centric captions and we provided human-collected confidence scores to measure the degree of commonsense expressed in the high-level captions. We also provided baseline results on generating captions for individual axes, as well as synthesised narrative captions by combining these three high-level axes of description.

Differently from current VL captioning datasets, the high-level captions capture the human interpretation of the scene allowing for inference and expectations. We discussed how they can be used also in combination with low-level captions to improve research in visual commonsense reasoning and multimodal grounding of visual concepts into linguistic expressions and for generative tasks. We will provide a practical study of this aspect in Chapter 4. We further exploit the alignment between high- and low-level captions in the HL dataset, in the explainability domain of VL models in generative settings, in Chapter 5. We also hope that the HL dataset would provide useful ground to foster future research in this direction.

## Ethical Considerations

The data collection received ethical approval from the University of Malta Research Ethics Committee, with reference number $7607 - 18012021$. This data is intended to be used for training, fine-tuning, and performing experimental evaluations of machine learning models. The dataset from which the images were originally sourced is a widely-studied, publicly available resource. As far as we are aware, the data does not contain harmful or offensive content. However, we acknowledge that any biases in the collection of images and/or captions in the original dataset will also be present in the HL Dataset.

# 4

# Analysing VL models' grounding capabilities

The material in this Chapter is based on:

- Michele Cafagna, Kees van Deemter, Albert Gatt. *What Vision-Language Models "See" when they See Scenes*, 2021, ArXiv preprint 2109.07301;

- Michele Cafagna, Kees van Deemter, Albert Gatt. *Understanding Cross-modal Interactions in V&L Models that Generate Scene Descriptions*. Proceedings of the Workshop on Unimodal and Multimodal Induction of Linguistic Structures, The 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP2022);

**Contributions**: Michele Cafagna: implementing and running the experiments; writing and revising the papers. Albert Gatt and Kees van Deemter: supervising the research, writing, and revising the papers.

# 4.1 | Introduction

The HL dataset introduced in Chapter 3 enables a direct alignment of object-level descriptive caption with high-level descriptions capturing the human interpretation of the image along three axes: *actions, rationales* and *scenes*. As briefly discussed in Section 3.7, this distinction allows a direct comparison of grounding capabilities of VL models when exposed to descriptions at different level of abstraction.

In this context, we delved into the *scene* axis of our resource, exploring how VL models ground scene descriptions. Specifically, we analyzed how these models relate to the objects within scenes and their visual arrangement. This investigation aims to shed light on how VL models ground scene descriptions in relations to the objects they contain and their visual layout.

Current research on human perception suggests that humans do not perceive scenes exclusively in terms of the objects they contain, and that visual salience is not exclusively determined by bottom-up features such as colour and texture (Malcolm et al., 2016). Rather, visual stimuli are considered "scenes" because their elements constitute a meaningful whole, both in terms of their contents (e.g. one expects an oven in a kitchen, but not in a living room) and in terms of their spatial arrangement (e.g. ovens do not typically hang from the ceiling).

These observations have provided the impetus to work showing that violations of scene "semantics" (content) and "syntax" (spatial arrangement) exact a cognitive cost during perception (e.g. Biederman et al., 1982; Võ and Wolfe, 2013; Võ, 2021). For instance, a typical scene of a dining room includes chairs around a table. Thus, a chair hanging off the ceiling would be a syntactic violation. Such violations lead to higher processing load in humans.

A related strand of modeling research in computer vision has shown that scene-level priors generate expectations about objects and their configurations, impacting the salience of objects in a way that classical, feature-based models of attention (e.g. Itti and Koch, 2001) would not predict (Oliva and Torralba, 2007; Torralba et al., 2006). Indeed, the problem of linking low-level features with high-level semantic information is an instance of the problem referred to as the "semantic gap" in computer vision (Ma et al., 2010).

Moreover, scene-level captions are less redundant with respect to the image they describe, but convey enough information to generate inferences about content and structure as they rely on implicit world knowledge. (For example, in a baseball field we typically expect to find players dressed in uniforms, placed in a certain way and spectators located on grandstands.) Thus, understanding if models succeed in linking scene-level

66

 **LN:** This is the picture of a stadium. In the foreground, there is a person [. . . ] At the back there are a group of people sitting [. . . ].
**COCO:** a baseball player getting ready to swing at a baseball game in a stadium packed with people.
**HL-scenes-1k:** The picture is shot in a baseball field

Figure 4.1: An example of a scene with COCO and Localized Narrative (LN) object-level captions, versus HL scene-level description.

and object-level descriptions helps to shed light on whether these models learned some of the relevant inferential links.

In Section 2.3.2, we emphasized that data used for VL pre-training usually contains highly descriptive text that mentions objects and their spatial relationships. For instance, the COCO (Chen et al., 2015c) and Localized Narratives (LN; Pont-Tuset et al., 2019) captions for Figure 4.1 are of this type, though they differ stylistically. By contrast, the third caption in the figure, from the scene-axis of the HL dataset introduced in Chapter 3, is what we refer to as "scene-level", focusing on *what type* of scene or location is depicted.

In Figure 4.1 both the object- and scene-level descriptions describe the picture, albeit in different ways. Indeed, it would be expected that, for a VL model to display true grounding capabilities, it should be able to match both types of descriptions with the image. For models that do display this capability, a natural follow-up question is whether their representations captures interesting connections between scenes on the one hand and the objects within them on the other.

In this Chapter, we investigate the capability of VL models to handle object-level and scene-level descriptions equally well. Positive evidence would suggest that such models are learning useful associations between the elements of a scene and the overall

scene type, as captured in textual descriptions. The Chapter is composed of two main Sections:

In Section 4.3, we perform an in-depth analysis in a zero-shot setting on three state-of-the-art pre-trained VL models with the goal of understanding what VL models learn, as a function of the data they are pre-trained on and the model architecture.

In Section 4.4, we further proceed with our investigation, by trying to shed light on the ability of VL models to reason about the relationship between scenes and their components when directly exposed to such descriptions in generative settings.

## 4.2 | Related Work

**Methods**   VL models have been extensively evaluated on tasks such as Visual Question Answering (Goyal et al., 2017) or image retrieval (Lin et al., 2014b). More recently, there has been increased interest in understanding the nature of the representations and capabilities learned by large, pre-trained models, for example via probe tasks or investigation of their attention heads (see Belinkov and Glass, 2019, for a survey). This has also been done for VL models. For example, Li et al. (2020b) consider VisualBERT's attention heads in a manner similar to Clark et al. (2019a), showing that it is able to ground entities and syntactic relations (see also Dahlgren Lindström et al., 2020; Ilharco et al., 2020). Hendricks and Nematzadeh (2021a) similarly seek to obtain an in-depth understanding of the representations learned by VL models, finding that they have difficulty with grounding verbs in visual data, compared to other morphosyntactic categories.

The present Chapter has a similar motivation but focuses on models' ability to reason in a grounded way about the relationship between entities and scenes. We focus on a number of methods: in Section 4.3 we develop an ablation method in the textual and visual modalities to uncover asymmetries in the extent to which VL models rely on textual or visual modalities, similarly to Frank et al. (2021a), which was developed concurrently with ours.

In Section 4.4 we focus on three techniques for model analysis: attention analysis, multimodal ablation and probing. Analyses of attention in pre-trained VLmodels include both quantitative methods (e.g. Abnar and Zuidema, 2020) and qualitative analysis (e.g. Li et al., 2020c; Wei et al., 2021). We use both methods to study how generative VL models deploy attention during the generation, of object-centric, versus scene-level captions.

More generally, a number of tasks have been developed to test the ability of VL models to reason with a combination of linguistic and visual cues, including VCR (Zellers

et al., 2019), SWAG (Zellers et al., 2018) and NLVR (Suhr et al., 2017a, 2019) which focus on visual common sense reasoning. VALSE (Parcalabescu et al., 2022b) tests specific visually grounded linguistic capabilities such as counting, spatial relations and coreference. The Winoground (Thrush et al., 2022) benchmark focuses on the visio-linguistic compositionality. Pezzelle et al. (2020a), in work complementary to our own, address the relationship between visual and textual modalities, exploring a task in which the text does not provide an object-level description of an image.

**Models**   Scene recognition is a central task in computer vision, with extensive work on scene categorisation systems (e.g. Anderson et al., 2021) and several datasets in addition to the ones used in this work, including ImageNet (Deng et al., 2009), Places (Zhou et al., 2014) and SUN (Xiao et al., 2010). However, there has been little work at the VL interface, exploring the capabilities of models to link scene- and object-level representations. Some precedents for the concerns addressed in this Chapter are found in the image captioning literature. For example, an influential proposal by (Anderson et al., 2018) combines top-down and bottom-up attention to combining local and global features. CapWAP (Fisch et al., 2020a) conditions image captioning on questions that determine which information is relevant to current communicative needs, going beyond object-level description. Closer to the scope of the work presented here, a recent pre-trained VL model, SemVLP (Li et al., 2021a), combines single- and dual-streams for feature-level and high-level semantic alignment.

However, recent studies have shown that such architectural differences in VL models, i.e. single- and dual-stream, lead to roughly the same performance under the same training settings (Bugliarello et al., 2021b). This is in line with the results of our analysis, performed in Section 4.3. We test three SOTA VL models at the time of writing this thesis, with different architectures and sizes, finding that training data and training objectives are relevant factors to achieve good scene grounding rather than architectures.

From the perspective of caption generation, the Oscar (Li et al., 2020e) single-stream architecture has emerged as an influential model. Oscar enforces grounding between image-caption pairs by using object labels as anchor points (a strategy also adopted by Hu et al., 2021a). This makes it particularly suited to the goals of the work presented in Section 4.4, namely, in-depth analysis of the cross-modal interactions in the treatment of objects during generation. Oscar and its successors, VinVL (Zhang et al., 2021) and LEMON (Hu et al., 2022a) achieved SOTA performance on captioning tasks such as COCO and `nocaps`.

| | Training size (# image-sentence pairs) | Model size (# parameters) | Pretraining Objectives |
|---|---|---|---|
| CLIP | 400M | 151M | ISA |
| VisualBERT | 330k | 112M | ISA, MLM |
| LXMert | 9.18M | 228M | ISA, MLM MOP, VQA |

Table 4.1: Comparison of training settings for the three models (**ISA**: Image-Sentence Alignment, **MLM**: Masked Language Modeling, **MOP**: Masked Object Prediction, **VQA**: Visual Question Answering)

# 4.3 | How do VL models "see" scenes?

In this Section, we investigate whether VL models are able to handle object-level and scene-level descriptions equally well. We perform an analysis in a zero-shot setting on three state-of-the-art pre-trained VL models. To our knowledge, this is the first systematic comparison of model capabilities on object- versus scene-level grounding.

The goal of this study is therefore not to establish new SOTA results, but to further our understanding of what VL models learn, as a function of the data they are pre-trained on and the model architecture. We chose three models differing in many settings (including training set size, architecture, number of parameters, and model size). All of the models were however optimized on the image-sentence alignment task.

## 4.3.1 | Models

Many VL models typically combine textual and visual features in a single or dual-stream architecture. Though the two architectures have been found to perform roughly at par when trained on the same data in comparable settings Bugliarello et al. (2020), in this work we include widely-used representatives of both at the time of writing, as we are interested in their zero-shot grounding capabilities in their original settings. We also include a third model which differs in structure and is trained on a much larger and more varied dataset. Table 4.1 gives an overview of some of the properties of the models we consider.

**LXMERT**   (Tan and Bansal, 2019) is a dual-stream model, which encodes text and visual features in parallel, combining them using cross-modal layers. LXMERT is trained on COCO captions (Chen et al., 2015c) as well as a variety of VQA datasets, with an image-text alignment objective, among others. We used the implementation of LXMERT

in the `transformers`[1] library.

**VisualBERT**   (Li et al., 2020c) is a single-stream, multimodal version of BERT (Devlin et al., 2019b), with a Transformer stack to encode image regions and linguistic features and align them via self-attention. It is pre-trained on COCO captions (Chen et al., 2015c). Image-text alignment is conceived as an extension of the next-sentence prediction task in unimodal BERT. Thus, VisualBERT expects an image $i$ and a correct caption $c_1$, together with a second caption $c_2$, with the goal of determining whether $c_2$ matches $\langle i, c_1 \rangle$. We use the publicly available implementation of the model.[2]

**CLIP**   (Radford et al., 2021a) combines a transformer encoder for text with an image encoder based either on Visual Transformer (Dosovitskiy et al., 2020a) or a Resnet network (He et al., 2015b) jointly trained using contrastive learning to maximise scores for aligned image-text pairs. CLIP is trained on around $400M$ pairs sourced from the Internet, a strategy similar to the web-scale training approach used for unimodal models such as GPT-3 (Brown et al., 2020b). We note that the visual backbone for this model differs from that of LXMERT and VisualBERT, both of which use Faster-RCNN (Ren et al., 2015a). In our setting we used the ViT-based visual backbone.

For all experiments, we truncated textual captions to a maximum length of 50 tokens, following standard practice for such models, including CLIP.

## 4.3.2 | Data

We used four different datasets for our experiments, which overlap to different degrees with the data that LXMERT and VisualBERT were trained on.[3] The extent of overlap is shown in Table 4.2. We used four datasets: two aligning descriptive captions (Localized Narratives, COCO) and two aligning scene descriptions to images (HL-scenes-1k, ADE20k). We used two different datasets per kind of caption (descriptive- vs scene-level) to account also for stylistic differences within the same kind of captions.

**Localized Narratives**   Localized Narratives (LN) Pont-Tuset et al. (2019) is a VL dataset created by transcribing speech from annotators who were instructed to give object-by-object descriptions as they moved a mouse over image regions. LN captions tend to be highly detailed and stylistically similar to speech. We used LN as a source of object-level

---

[1]`github.com/huggingface/transformers`

[2]`https://github.com/uclanlp/visualbert`

[3]CLIP was trained on a web-harvested dataset not released to the public, thus its composition is unknown.

|              | LXMERT | | VisualBERT | | CLIP | |
|--------------|:---:|:---:|:---:|:---:|:---:|:---:|
|              | I | C | I | C | I | C |
| COCO         | ✓ | ✓ | ✓ | ✓ | ? | ? |
| LN           | ✓ | ✗ | ✓ | ✗ | ? | ? |
| HL-scenes-1k | ✓ | ✗ | ✓ | ✗ | ? | ✗ |
| ADE20K       | ✗ | ✗ | ✗ | ✗ | ? | ? |

Table 4.2: Presence/absence (<✓>/<✗>) of the (I)mages and (C)aptions of the datasets used for the experiments in the training data of VisualBERT and LXMERT. Although the composition of CLIP's training data is unknown (marked with <?>), the HL-scenes-1k captions are certainly not included in it, as they were collected after CLIP's release.

captions. The images in LN come from pre-existing datasets; this allows us to align LN captions with images and captions from datasets such as COCO and ADE20K.

**COCO**   (Lin et al., 2014c) consists of images paired with captions and object annotations. LN captions are also available for the same images. We used images and captions from the 2017 COCO validation split, as well as the corresponding LN captions.

**HL-scenes-1k**   **H**igh **L**evel scenes is a subset of the *scene* axis of the HL dataset (Cafagna et al., 2023b), introduced in Chapter 3. HL-scenes-1k is composed of 1k images, each depicting at least one person, sampled from the 2014 COCO `train` split. As described in Section 3.3 we crowd-sourced three annotations per image on Amazon Mechanical Turk, showing crowd workers the image and asking them to write a description in response to the question *Where is the picture taken?* Crowd workers were asked to respond using full sentences and it was made clear to them that their answer to this question should bring to bear their knowledge of typical, or common, scenes. Figure 4.2 shows an image with three different scene descriptions. The scene-level captions are then aligned with the original COCO captions.

**ADE20K**   (Zhou et al., 2017) is a computer vision dataset containing 20k images comprehensively annotated with objects, parts and scene labels. We used ADE20K as a source of scene-level captions. For our experiments, we filtered out images with scenes that in the dataset are labeled as `unknown`. We produced captions for each image using a simple template-based generation method, whereby a scene label is inserted into one of the templates below:

- *it is a* SCENE

72

- *this is a* SCENE

- *it is located in* SCENE

We aligned the resulting scene-level descriptions and the corresponding ADE20K images to the corresponding object-level captions in LN.



**Where is the picture taken?**
- in a bedroom
- the picture is taken in a bedroom
- this is the bedroom

Figure 4.2: COCO image with three HL-scenes-1k scene descriptions.

|              | # images | caption source | # captions per source |
|--------------|----------|----------------|-----------------------|
| HL-scenes-1k | 1000     | HL<br>COCO     | 1000<br>1000          |
| ADE20k       | 19733    | ADE20K<br>LN   | 19733<br>19733        |
| COCO         | 5000     | COCO<br>LN     | 5000<br>5000          |

Table 4.3: Statistics for the HL-scenes-1k, ADE20k and COCO.

We corrected the HL-scenes descriptions for possible typos using the Neuspell Toolkit (Jayanthi et al., 2020). Finally, we paired our scene-level HL-scenes-1k captions with the previously available COCO and LN object-level captions. Figure 4.1 provides an example.

|        |             | LXMERT | CLIP | VisualBERT |
|--------|-------------|--------|------|------------|
| Object | ADE20k + LN | 28.4   | **96.8** | 39.0   |
|        | COCO + LN   | 59.1   | **98.7** | 65.2   |
|        | COCO Cap.   | 79.3   | **99.1** | 64.4   |
| Scene  | ADE20k      | 58.0   | **97.6** | 17.3   |
|        | HL-scenes-1k | 45.5  | **91.5** | 55.3   |

Table 4.4: Image-sentence alignment accuracy on object-level and scene-level captions. Chance performance is at 50% (*LN* = Localized Narratives).

Dataset statistics are shown in Table 4.3. For ADE20k, the numbers are for images that are not labeled as having an *unknown* scene. In COCO, there are five captions associated with each image.

### 4.3.3 | Image-sentence alignment experiments

We first tested models in the image-sentence alignment task on both object- and scene-level descriptions. Since we are interested in the capabilities of the pre-trained models, and since pre-training included alignment for all models (see Table 4.1), we did not fine-tune them. Rather, we used the models' pre-trained alignment head to predict whether a scene-level or object-level caption correctly describes an image, or not; negative samples were randomly drawn. [4] See Appendix B.1 for full details.

Table 4.4 shows that LXMERT and VisualBERT perform adequately on object-level COCO Captions, though performance is lower than would be expected, given that they were pre-trained on this dataset. In the case of LXMERT, one possible explanation is catastrophic forgetting, arising from the fact that this model is pre-trained for its final ten epochs on VQA (similar observations are made by Parcabalescu et al., 2021). For both models, performance drops dramatically on LN captions. This is likely due to a stylistic difference: compared to COCO captions, LN captions are longer, more discursive, and contain disfluencies.

In contrast, CLIP performs close to ceiling on all three datasets, possibly reflecting the benefits accrued from the size and diversity of its pre-training data.

On scene-level captions, performance is somewhat above chance for LXMERT on ADE20k template-based descriptions, and for VisualBERT on HL-scenes-1k. Otherwise, performance is below 50% for both models. Once again, CLIP performs above 90%, though there is a drop in performance from the template-based ADE20k descriptions to

---

[4] Note that this setting is the same used by the models in their pre-training.

human-authored HL-scenes-1k scene-level captions, possibly reflecting the more predictable nature of the former.

### 4.3.4 | Ablation experiments on CLIP

CLIP is the only one of the three models which is successful at matching scene-level and object-level captions to images. Motivated by findings in Section 3.4.1, showing that scenes tend to be correlated with the presence of certain objects, we probed CLIP's capabilities further, paying particular attention to the question whether CLIP links *scene types* (e.g. kitchen) to *scene contents* (e.g. oven, pizza) in image-text matching.

Whereas a standard image-text alignment setup compares the model's success at aligning actual versus random captions with images, here we directly compared the preference of the model for scene- versus object-level descriptions, as a function of (i) the entities mentioned in the object-level caption; (ii) the entities visible in an image. To this end, we used textual and visual ablation on captions and images; an example is shown in Figure 4.3.

**Textual ablation**   Given an object-level caption, we identified all the Noun Phrase (NP) in the caption and create new versions by removing each possible subset of the set of NPs, with the restriction that the resulting caption must always contain at least one NP. When NP removal resulted in dangling predicates, we removed them to preserve grammaticality. For example, in Figure 4.3, when the NP *"A man"* is ablated from the caption, the dangling predicate *"rides"* was also removed. NPs were detected with Spacy v.3, using the pipeline for English with the `en_core_web_md` pretrained models. The right panel of Figure 4.3 shows the original caption and examples of ablated captions.

For a given image $i$ with object-level caption $o$ and scene-level caption $s$, we compared how $P(o|i)$ – CLIP's estimate of the probability that $o$ matched $i$ – changes as NPs were removed from $o$, and to what extent this caused CLIP to assign higher probability $P(s|i)$, to $s$ as the match for $i$. We report two comparisons, one on LN captions versus ADE20K template-based scene descriptions; and one on COCO captions against HL-scenes-1k scene-level descriptions.

To control for possible loss of grammaticality after ablation, we scored ablated captions with GRUEN (Zhu and Bhat, 2020), a BERT-based model which has been shown to yield scores that correlate highly with human judgments.[5] CLIP probabilities for ab-

---

[5]GRUEN returns a combined score consisting of a linear combinaton of grammaticality, focus and Coherence. Here, we used only the grammaticality scores.

**COCO:** A man rides a motorcycle on a road through a grassy, hilly area.

**Ablated Captions:**

- a grassy, hilly area *(A man, a motorcycle, a road)*
- a road *(A man, a motorcycle, a grassy, hilly area)*
- a road through a grassy, hilly area *(A man, a motorcycle)*
- A man rides a road *(a motorcycle, a grassy, hilly area)*
- a motorcycle a grassy, hilly area *(A man, a road)*
- A man rides a grassy, hilly area *(a motorcycle, a road)*
- A man rides a motorcycle *(a road, a grassy, hilly area)*
- A man rides a road through a grassy, hilly area *(a motorcycle)*
- a motorcycle on a road *(A man, a grassy, hilly area)*
- A man rides a motorcycle on a road *(a grassy, hilly area)*
- A man rides a motorcycle a grassy, hilly area *(a road)*
- a motorcycle on a road through a grassy, hilly area. *(A man)*

Original Image

Occluded Image

Figure 4.3: Example of visual and textual ablation. *Left*: Original image and image with occluded object. *Right*: Original caption and different ablated captions. NPs removed are shown in parentheses.

lated textual captions yielded a significant, but very low correlation with grammaticality (Pearson's $r = 0.1, p < .01$) suggesting that grammaticality did not affect the scores.

**Visual ablation**    Given an object-level caption and an image, we extracted all nouns from the caption and extract the embedding vector for each noun using pretrained Fast-Text embeddings.[6] We passed the image through the Faster-RCNN object detector[7] to detect entities. We extracted embeddings for each entity label. Then, we identified regions to be masked by comparing embeddings for entity labels $l_e$ against embeddings for nouns $n_e$ in the caption, considering them a match if $\text{cosine}(l_e, n_e) \geq 0.7$. This thresh-

---

[6]We used the model with $2M$ word vectors trained with sub-word information from Common Crawl `https://fasttext.cc/docs/en/english-vectors.html`

[7]Faster R-CNN ResNet-50 FPN pre-trained on COCO, available from the `torchvision` module in `Pytorch`

old was empirically determined by maximizing the number of correct matches in a representative sample of data. Bounding box regions corresponding to matched entities were occluded with a greyscale mask. The left panel of Figure 4.3 compares the original and masked image.

Once again, we are interested in whether CLIP's estimate of the alignment probability of object- versus scene-level captions, changes as elements of the visual input are masked.

|     | ADE20k | HL-scenes-1k |
|-----|--------|--------------|
| T   | 205k   | 10027        |
| V   | 10788  | 625          |
| V+T | 1078   | 625          |

Table 4.5: Total number of ablations generated per dataset, across all the ablations experiments using T(extual) ablation, V(isual) ablation, and both (V+T).

Table 4.5 provides the number of ablations analysed in the study. For both ADE20k and HL-scenes-1k, we obtained a number of ablated captions that is greater than the respective dataset sizes in Table 4.3 because for each example we generated all the possible combinations of noun phrases. For the Visual and Visual+Textual ablations, the number of ablated instances is lower than the dataset size, because we omitted all the images where no object is detected.

**Results** The results of image-sentence alignment using CLIP, after ablation, are shown in Table 4.6. With no ablation, the model assigns a higher probability to object-level descriptions, suggesting that CLIP has higher confidence in aligning an image-text pair when the text focuses on objects rather than scenes. This preference is far more marked for COCO/HL-scenes-1k, in line with the observation (Table 4.4) that HL-scenes-1k scene descriptions are somewhat more challenging for this model.

As entity-level information is removed from the object-level caption (row T in Table 4.6), the model assigns higher probability to the scene-level caption, suggesting that the model leverages the visual information to align with the scene description.

In contrast, visual ablation (row V) results in the opposite tendency: when entities are occluded in the image, the model assigns a higher probability to object-level captions compared to scene-level descriptions.

These results suggest that CLIP aligns images to scene-level descriptions based on the entities visible in the images. As these are masked in the image, entity-level captions

|             | ADE20k | | HL-scenes-1k | |
|-------------|--------|-------|------|-------|
|             | LN     | Scene | COCO | Scene |
| No ablation | 55.6   | 44.4  | 95.7 | 4.3   |
| T           | 22.0   | 78.0  | 67.2 | 32.8  |
| V           | 74.9   | 25.1  | 71.2 | 28.8  |
| V+T         | 68.4   | 31.6  | 63.3 | 36.7  |

Table 4.6: CLIP preferences for object-level versus scene-level captions for two different image datasets (i.e. ADE20k and HL-scenes-1k). Within each sub-column we show the model's preference rate when the same image is paired with a scene-level (e.g. Scene) vs object-level (e.g. LN) caption. Each row corresponds to different setups, namely when performing no ablation, T(extual) ablation, V(isual) ablation, or both (V+T). Note that each cell sums-up to 100.



(a) Scene: *kitchen*

(b) Scene: *road*

(c) Scene: *room*

(d) Scene: *park*

Figure 4.4: Visualisations of entities (*e*) in four different scene types (*s*). Font size is proportional to $P(s|e)$

are aligned with higher probability. On the other hand, when both sources of information are ablated, CLIP once again assigns a higher probability to object-level captions.

## 4.3.5 | Scenes vs. entities

Our findings suggest that CLIP reasons about scenes on the basis of salient objects within them. If this is the case, then the probability assigned by clip to an image-scene caption pair should diminish as more salient entities are visually ablated in the image. This is also motivated by the correlation found between scenes and diagnostic objects present in the image, observed in the experiments in Section 3.4.1.

To investigate this further, we used scene labels extracted from the HL-scenes-1k captions and the object detections produced for the visual ablation (Section 4.3.4). For a scene label $s$ and entity label $e$, we computed $P(s|e)$ as follows. Let $e$ be an entity detected $n_e$ times in the dataset, of which $n_{e,s}$ times in images depicting scene $s$. We computed:

$$P(s|e) = \frac{n_{e,s}}{n_e}$$

Figure 4.4 shows visualizations for entities detected in four example scene types found in the HL-scenes-1k dataset.

For all images with at least three detected entities, we considered the image-sentence alignment probability assigned by CLIP to the *scene*-level description, when the top 1, 2, or 3 most likely entities in the scene are masked. We therefore averaged over those images containing at least three detected entities (53/174 total scenes).

Figure 4.5 displays the average alignment probability assigned by CLIP to images and scene-level captions, as entities are progressively masked in the image. The figure displays a linear trend, with the probability dropping as more likely entities are removed. A one-way Analysis of variance (ANOVA) comparing the change in log probability as 1, 2, or 3 entities are removed showed that the difference is significant ($F(2, 156) = 4.25, p < 0.05$).

Thus, when CLIP aligns images with scenes, it relies on object-level information in the visual modality. This explains why the removal of object mentions in text results in a higher preference for scene-level descriptions since the objects are detectable in the image. By the same token, masking objects in images causes the model to rely more on the entity-level information in the text.

## 4.3.6 | Effect of length and informativeness

So far, our analysis suggests that CLIP reasons about scenes based on object-level information. However, the length of the caption might be a possible confounding factor. Some of our results might simply be due to the model assigning a higher alignment

Figure 4.5: CLIP scene-level description probabilities after masking top 1-3 entities. Error bars represent standard deviations.

probability to a caption that is longer or more informative. This could provide an alternative explanation for the changes observed above in the alignment probabilities after textual ablation.

To account for this, we replicated the alignment experiment using single words. Once again, we used the scene labels extracted from HL-scenes-1k scene-level descriptions and identify the top three most likely entities in a given scene, as in the previous experiment (see Figure 4.5).

Given an image, we compared image-text alignment probabilities in CLIP for single-word object labels (e.g. *motorbike*) and single-word scene labels (e.g. *road*). In this setting, CLIP displays a moderate preference for scene labels (63%), suggesting that such labels are more informative than object-level labels, for the one-word alignment task.

We performed a qualitative analysis, inspecting 5 cases where the model has a clear preference for scene label or object label. Some examples are shown in Figure 4.6. CLIP assigns a higher probability to object labels when images have salient, foregrounded entities. When entities are less salient or in the background, the model prefers scene labels. See Appendix B.2 for more examples.

*resort*: 3% *person*: 97%                    *resort*: 99% *snowboard*: 1%

*kitchen*: 99% *bowl*: 1%                    *kitchen*: 11% *knife*: 89%

Figure 4.6: Scene vs entity one-to-one comparison. In the top left image, there are many people in the foreground and the entity *person* is preferred over the scene label *resort*. At the top right image, people are snowboarding in the background and the scene label is preferred over the entity label *snowboard*. Similarly at the bottom, when the *knife* is in the foreground (bottom-right), it is preferred over the scene label *kitchen*. When the object is less evident, such as the *bowl* (bottom-left), the scene-label is preferred.

## 4.3.7 | Conclusions

In order to address the symbol grounding problem, VL models should be able to capture the relationship between an "object-level" view of an image, focusing on objects and their configuration, and the higher-level scene it corresponds to. In this study we found that when models do this, they rely on object-level information in the *visual* modality, to link images to scene descriptions in the *textual* modality; this is influenced by the probability of entities occurring in particular scene types.

Of the models tested, we find that LXMERT and VisualBERT perform poorly on this task, and also suffer when captions deviate stylistically from their pre-training data. For these models, testing on ADE20k, amounts to a full zero-shot setting, whereas for Localized Narratives and HL-scenes-1k, this only applies to the textual input, as the

images are included in their training data. With the exception of HL-scenes-1k, a new dataset, it is an open question, whether testing for CLIP was zero-shot since this model was trained on web-scale data, which is often unfathomable (Bender et al., 2021). On the other hand, model size is clearly not the determining factor; CLIP has fewer parameters than LXMERT, for example (cf. Table 4.1).

We believe that two additional factors contribute to the success of CLIP. First, its contrastive learning objective may result in greater sensitivity to fine-grained distinctions between captions for image-sentence alignment. A second feature is its visual backbone, which (in the version used in our experiments) is based on Vision Transformer (ViT Dosovitskiy et al., 2020a). Recently, BERT-inspired architectures have achieved notable success on computer vision tasks (see also Bao et al., 2021b). Tuli et al. (2021) have shown that ViT is more consistent with characteristics of human vision than a convolutional network, extracting image features that are not strictly local. This could partially underlie the model's ability to use visual object-level information to align with scene-level captions.

## 4.4 | Cross-modal relationships in scene descriptions

While in the previous Section, we investigate the models' capability to create meaningful object-scene relationships, here, we focus on generative VL model. In Chapter 3 we showed that it is possible to fine-tune image captioning models to the different axes. In this section, we delve into the underlying mechanisms within these models when such fine-tuning is applied.

We present a study of object-centric versus scene-level captioning, focusing on the impact of the exposure of pre-trained VL models to scene-level descriptions. We focus on VinVL (Zhang et al., 2021), a BERT-based model in the OSCAR family (Li et al., 2020e) of models, which have recently dominated the state of the art in image captioning.[8]. As already shown in Section 3.5, VL models trained on object-centric captions can easily adapt to scene-level descriptions. In this study, we go beyond metric-based results, diving into the mechanisms driving the model towards generating scene-level descriptions. Moreover, we ask whether insights yielded by this analysis are compatible with the findings reported in Section 4.3.

The main contributions of this Section are:

---

[8]At the time of this work, three OSCAR-based models (OSCAR, VinVL, LEMON) were among the top 5 in the leaderboard of the COCO image captioning task.

i) We performed an in-depth investigation of the impact of fine-tuning on the pre-trained model. The analysis is designed to thoroughly inspect object-scene relations by exploiting cross-modal attention (Section 4.4.3), coupled with probing (Section 4.4.5) and ablation studies (Section 4.4.4).

ii) We show that (i) VinVL's pre-trained representations are rich enough to support scene-level captioning, but that (ii) fine-tuning results in a different deployment of attentional resources. This bears parallels to the findings in Section 4.3, where we show that scene descriptions understanding, relies on object-level information and research on human scene perception.

## 4.4.1 | Data

We used the HL-scenes dataset, namely the entire *scene* axis of the HL dataset introduced in Chapter 3. HL-scenes thus constitutes a superset of HL-scenes-1k presented in Section 4.3.2. It is composed of 14,997 image-caption pairs, split into 11,999 for training and 1,499 each for validation and testing as described in Chapter 3. This dataset is particularly suitable for the analysis presented here, as it aligns object-centric captions from COCO, with crowd-sourced scene-level descriptions. Moreover, differently from the study presented in Section 4.3, we are less interested in the stylistic variations within the same kind of caption. In fact, we will be focusing on the multimodal interplay between images of a scene vs object-centric captions.

## 4.4.2 | Model

VinVL (Zhang et al., 2021) is a single-stream BERT-based model with a Faster-RCNN (Ren et al., 2015b) visual backbone. It is an extension of Oscar (Li et al., 2020e). VinVL implements a training strategy where object tags are used as anchor points between the visual and textual modalities to facilitate cross-modal alignment. As pointed out by Li et al. (2020e), this strategy is motivated by the fact that in the datasets used to pre-train multimodal models, between 1 and 3 of the objects detected by the visual backbone are mentioned in the caption. However, the object labels are provided by an off-the-shelf object detector separately trained on Visual Genome (Krishna et al., 2017b). VinVL was pre-trained on a combination of COCO (Chen et al., 2015c), Conceptual Captions (Sharma et al., 2018a), SBU captions (Ordonez et al., 2011a), and Flickr30k (Young et al., 2014a), as well as additional VQA data.

**COCO**
*Reference:* a close-up of a kitten looking and a dog laying
in the background.
*Generated:* a cat and a dog sitting next to each other.

**HL-scenes**
*Reference:* in the home.
*Generated:* the picture is taken in a house.

Figure 4.7: Scene-level captions in HL-scenes, with corresponding object-centric COCO caption. The generated captions are outputs from VinVL before and after fine-tuning (see Section 4.4.2).

**Fine-tuning**   We first established that VinVL can generate scene descriptions after fine-tuning, before turning to an in-depth analysis of the model's attention and internal representations. We expected reasonably good results in view of the experiments conducted in Section 3.5 with baseline models.

We noted that since the HL-scenes dataset extends the COCO dataset, the model has been exposed to the images of the HL-scenes dataset during pre-training on COCO. On the other hand, the scene descriptions are completely novel. We fine-tuned on scene descriptions for 10 epochs. We used the standard configuration used by Zhang et al. (2021) for image captioning. At inference time, we fixed the maximum generation length to 20 tokens and use a beam size of 5.

We fine-tuned the VinVL pre-trained base version[9] using the original configuration for 10 epochs on scene descriptions. We refer to it as the *fine-tuned* model. Since the HL-scenes dataset images are included in COCO, we used the pre-computed visual features and labels provided in the original VinVL implementation. We refer to the *pre-trained* model, as the base model trained on the image captioning task on COCO captions optimized using cross-entropy. All the experiments involving the pre-trained

---

[9]`https://github.com/microsoft/Oscar/blob/master/VinVL_MODEL_ZOO.md#`
`Oscarplus-pretraining`

Figure 4.8: Inbound attention of the `[SEP]` per input type token across the layers. Special tokens correspond to `[CLS]`, `[PAD]` and `[SEP]`.

| Epoch. | Bleu-4 | METEOR | ROUGE-L | CIDEr | SPICE |
|--------|--------|--------|---------|-------|-------|
| 2 | 49.3 | 29.3 | 67.1 | 161.8 | 32.6 |
| 4 | 49.7 | 30.1 | 68.1 | 168.5 | 34.0 |
| 6 | 48.5 | 29.8 | 67.3 | 164.9 | 33.5 |
| 8 | 48.9 | 30.2 | 67.6 | 165.8 | 33.9 |
| 10 | 49.1 | 30.4 | 67.7 | 168.0 | 34.4 |

Table 4.7: Automatic metrics computed over different epochs on the HL-Scenes validation set.

model were performed using the original configuration used in Li et al. (2020e). The fine-tuning was carried out with batch size 32 on an NVIDIA GTX 2080 TI 11 GB.

VinVL shows a quick adaptation to the scene-level descriptions from the first epoch. This adaptability recalls observations made for other transformer-based generative models (e.g. Brown et al., 2020a). We show an example in Figure 4.7. For completeness, Table 4.7 reports the automatic evaluation metrics computed on the validation set over 10 epochs.

(a) Attention matrix of the pre-trained model          (b) Attention matrix of the fine-tuned model

Figure 4.9: Attention matrices comparison for the image in Figure 4.7. We highlight the sub-blocks corresponding to vision-to-vision, vision-to-label and label-to-vision. In the pre-trained model, attention mass is sharply focused on individual portions of the input; after fine-tuning, a more even distribution is observed.

## 4.4.3 | How does attention to objects change from object-centric to scene-level generation?

We first investigated the model's self-attention before and after fine-tuning on the scene-level caption generation task.

**Method**   We focused on the self-attention patterns in the first layer, as they are directly connected to the inputs and do not depend on higher-level interactions which might obscure the fundamental changes in attention across the two modalities (visual features and labels) in VinVL. A discussion of attention patterns at higher layers can be found in Appendix B.3. We selected 100 random samples from the HL-scenes test-set and extracted the attention matrices before and after fine-tuning on scene descriptions. We aggregated the attention values by taking the maximum across all the heads, as it allowed us to observe where the model tended to assign a significant amount of attention, giving us a better view of the potential impact of fine-tuning on scene-level captions. VinVL prevents textual inputs from directly interacting with the other modalities during generation; therefore there is no interaction between caption tokens and visual features. On the other hand, the model includes object tags as anchors and this allows us to study the multimodal interactions between the visual features and these object labels.

86

Figure 4.10: Kernel density estimate of distributions of standard deviations against attention mass for pre-trained and fine-tuned VinVL.

**VinVL acquires a holistic view of the scene after pre-training**  Figure 4.9 is a representative example of self-attention matrices extracted from the pre-trained (4.9a) and fine-tuned (4.9b) model with the image in Figure 4.7. The pre-trained model, which generates an object-centric caption, focuses attention on individual input tokens in the vision-to-vision, vision-to-label and label-to-vision sub-blocks. After fine-tuning, as the model generates a scene-level caption, the self-attention appears to be more evenly distributed over the inputs (4.9b). This suggests that when generating scene-level captions, the model leverages a wider range of visual features with less exclusive focus on individual objects or labels.

We performed a quantitative analysis of the self-attention in the sub-blocks of the matrix involving visual regions and object labels, computing a kernel density estimate of the distributions of the standard deviations against attention masses for each of the 100 samples, where the attention mass is computed summing all the attention scores in each location along the heads dimension. The result is shown in Figure 4.10. It is clear that the fine-tuned model has an overall lower standard deviation than the pre-trained model. This confirms that a similar attention mass is distributed more evenly after fine-tuning. We take this as evidence that in the process of generating scene descriptions, the fine-tuned model acquires a more holistic view of the input image, in contrast to the highly object-centered deployment of attentional resources evident in the pre-trained model.

**VinVL relies on diagnostic objects when generating scene-level captions** VinVL re-distributes self-attention over a wider range of visual features after fine-tuning. Never-theless, previous work on scene perception (Self et al., 2019; Vo, 2021) leads us to expect that in describing a scene, the model needs to rely on highly diagnostic objects. This is consistent to our findings discussed in Section 4.3. We computed diagnosticity empiri-cally, based on the occurrence of objects in scenes in our dataset.

As described in Section 4.3.2, during the data collection, the annotators were asked to answer the direct question: *Where is the picture taken?* As a consequence, the scene captions often have a regular structure, captured by the following three representative examples:

- the picture has been taken in a *restaurant*

- on a *beach*

- this is in an *airport*

Let $S$ be the set of the $k$ most frequent scene types mentioned in scene-level captions in the HL-scenes dataset.

To extract the scene labels, we tokenized the scene captions and we removed punctu-ation and stop-words (we added the word *picture* to the list of the standard stop-words). Among the remaining tokens, we extracted all the nouns and we reduced them to lem-mas, then we computed the frequencies of the remaining tokens. This allowed us to extract the scene types (*restaurant*, *beach* and *airport*) from the captions, such as those shown in the examples above. The whole procedure was performed using spaCy[10].

We proceeded as follows:

1. $\forall s \in S$ we built $O_M^s = [o_1^s, o_2^s, ..., o_n^s]$, the ranked list of the $n$ most attended objects by the model $M$ when generating a description of a scene of type $s$.

2. Similarly, $\forall s \quad S$ we collected $O_D^s = [o_1^s, o_2^s, ..., o_n^s]$, the ranked list of the most frequent objects in images depicting scenes of type $s$ in the dataset $D$.

We measured the overlap between $O_M^s$ and $O_D^s$ by computing their Intersection over Union (IoU), which is only sensitive to overlap in content, as well as their Rank Biased Overlap (RBO) (Webber et al., 2010)[11], which is a similarity metric for ranked lists.

---

[10]https://pypi.org/project/spacy/
[11]https://github.com/changyaochen/rbo

| Scene | RBO @ | | | IoU @ | | |
|---|---|---|---|---|---|---|
| | 3 | 5 | 7 | 3 | 5 | 7 |
| station | 0.88 | 0.84 | 0.87 | 0.5 | 0.66 | 1.0 |
| road | 1.0 | 0.9 | 0.91 | 1.0 | 0.66 | 1.0 |
| room | 0.27 | 0.25 | 0.24 | 0.2 | 0.11 | 0.18 |
| sea | 0.88 | 0.84 | 0.8 | 0.5 | 0.66 | 0.55 |
| resort | 0.72 | 0.7 | 0.7 | 0.5 | 0.42 | 0.55 |
| house | 0.38 | 0.5 | 0.53 | 0.5 | 0.42 | 0.55 |
| restaurant | 0.55 | 0.55 | 0.54 | 0.5 | 0.42 | 0.53 |

Table 4.8: Rank Biased Overlap (RBO) and Intersection over Union (IoU) of the most attended objects and the most frequent objects for the top seven common scenes. Both metrics range from 0 (no overlap) to 1 (perfect correspondence).

RBO computes the similarity of two ranked lists, as follows:

$$RBO(S, T, p) = (1 - p) \sum p^{d-1} A_d \tag{4.1}$$

where $d$ is the depth of the ranking being examined, $A_d$ is the agreement between $S$ and $T$ given by the proportion of the size of the overlap up to $d$, and $p$ determines the contribution of the top $d$ ranks to the final RBO measure. We used the standard value of $p = 1$.

Table 4.8 shows RBO and IoU for the top 3, 5 and 7 objects in the lists. We observe that the two metrics correlate strongly ($r(19) = .81, p < .001$). From this, we conclude that during the generation of scene-level captions, the model attends more to diagnostic objects, i.e. those that are common in a scene of a given type. Moreover, we observe high scores for scene types such as *station, road, resort, sea*. In our dataset, these are characterised by frequently occurring objects, which are therefore highly diagnostic of scene type. In contrast, for scenes like *room, house, restaurant* we observe lower scores. We hypothesise that this is due to the fact that such scenes can contain a wider variety of objects, which individually have lower diagnosticity with respect to the scene type.

### 4.4.4 | How reliant is the model on diagnostic objects?

The results from the previous sections established that, following fine-tuning on scene-level descriptions, VinVL distributes attention more evenly over objects in a scene. Nevertheless, the objects that are most likely to be present in a scene attract the highest proportion of the attention mass. Our findings in Section 4.3 show that robust VL encoders like CLIP can be significantly affected by the ablation of diagnostic objects from the image when aligning scene-descriptions. This raises the question whether, by removing

| Scene | Top informative objects |
|---|---|
| restaurant | french fries, fork, submarine sandwich |
| road | vehicle number plate, traffic sign, traffic light |
| sea | surfboard, watercraft, boat |
| room | computer mouse, nightstand, tablet computer |
| station | train, suitcase, luggage and bags |

Table 4.9:  Most informative objects for some scenes ranked using PMI.

highly diagnostic objects from an image, the model representations are still informative enough to detect what type of scene is represented in an image in the generative setting.

We first address this issue from the perspective of generation: does a model fine-tuned on scene descriptions still manage to correctly describe a picture at the scene level, when highly diagnostic objects are unavailable? Given the more even distribution of attention observed across scene components in the fine-tuned model, our hypothesis would be that even in the absence of such highly diagnostic objects, the model can rely on other information to detect the scene type. Hence, we expect the fine-tuned model to be more robust to object ablation in the visual modality, compared to the model pre-trained on object-level captions.

As explained in Section 4.4.2, in VinVL, two separate models are used to (i) extract visual features corresponding to regions via the model's visual backbone; and (ii) to determine the object labels that function as anchors between the visual and textual modalities. This means we do not have an exact correspondence between object labels and visual features.

**Visual feature tagging**   For simplicity we refer to $vf$ as the bounding box a visual feature corresponds to, and $ot$ as the bounding box an object label corresponds to. To perform an ablation, we first established an approximate correspondence between $ot$ and $vf$, using $ot$ as a reference to assign an object label to the visual features.

We computed the IoU[12] between $vf$ and $ot$ and empirically assigned a label to a visual feature if $IoU(vf, ot) >= 0.6$. This threshold was determined as a trade-off between correct matches and noise, by manually testing on a sample of data. Moreover, if $vf$ is contained by or overlaps with $ot$ by at least 80% of its area, we assigned to $vf$ the label of $ot$. With this heuristic, we covered 74% of the visual features of every image of our sample.

---

[12]Note that in this section we refer to the Intersection Over Union to compute the overlap between two bounding boxes, not the metric used to compute the overlap between two sets of items as done in Section 4.4.3.

| # Ablation | Train-Val | Test |
|:---:|:---:|:---:|
| no ablation | 13498 | 1499 |
| 1 | 4269 | 469 |
| 2 | 2565 | 274 |
| 3 | 1554 | 170 |

Table 4.10: Sample size of the Train-Val and Test split after ablation of the top 1,2 and 3 most informative objects in the most frequent scenes. The top row corresponds to the original dataset split sizes.

**Computing object diagnosticity**   We used the scene labels extracted from captions in Section 4.4.3, and computed the PMI between scene types and object labels, similarly to what performed in Section 4.3 and Chapter 3. Examples of the most informative objects for some scenes are shown in Table 4.9.

**Ablation**   The ablation of an object was performed by removing its corresponding label from the list of object tags, which was replaced by a [PAD] token. All the visual features assigned to that object were removed by setting to 0 all their vector representations, similarly to (Frank et al., 2021b). We compared captions generated by both the pre-trained and fine-tuned model with and without ablation of the top 1, 2, and 3 most informative objects for a given scene in the test-set. As a result, an image was included in the ablation study if (i) it belonged to the set of most frequent scenes; and (ii) it contained the objects we wanted to ablate. This means that the higher the number of objects ablated, the smaller the sample of images matching these constraints. As shown in Table 4.10, the number of images matching this constraint is reduced up to 170 when 3 objects are ablated.

**Results**   We expected to observe some differences in the generations when ablation is applied, especially in the pre-trained model, as the ablation removes information that is explicitly verbalised in object-centric captions. In order to have a measure of the change we directly compared the captions generated by the pre-trained and the fine-tuned model, counting a change when the strings are different. For the pre-trained model, object-centric captions change 41% of the time after ablation, compared to 13% of the time for the scene-level captions by the fine-tuned model.

A manual inspection of a sample of items suggested that the changes in the captions involve minimal semantic shifts, often due to minor function word changes or a more generic term being generated for the noun denoting the scene type. Some examples are shown in Figure 4.11.

the picture is shot in a ski resort → the picture is taken in a snowfield *(jacket, tree, footwear)*
the picture is shot in a baseball field → the picture is taken in a ground *(sports uniform, man, boy)*
in a kitchen → in the kitchen *(kitchen appliance, countertop, cabinetry)*

Figure 4.11: Changes to scene-level captions generated by the fine-tuned model after ablation of three diagnostic objects. Ablated objects are shown in parentheses.



Figure 4.12: Confidence scores of the unchanged caption after ablation. On the left, the model generating scene-level descriptions (fine-tuned); on the right, the model generating objective descriptions (pre-trained).

In summary, the model is resilient to ablation in the visual modality, suggesting that its representations are robust for both types of generation tasks, but more so for scene-level captioning. This confirms the hypothesis based on findings of the attention analysis reported in Section 4.4.3, namely that when generating scene descriptions the model relies on a greater number of objects with lower individual diagnosticity for the scene type. This results in a higher resilience to object ablation in the visual modality.

We study the robustness of representations in more detail using probes, in Section 4.4.5.

Figure 4.13: Confidence shift of the unchanged captions when ablating the top 1, 2, and 3 most informative objects from the scene. A negative shift means that the caption was generated with higher confidence after ablation. On the left, the model generating scene-descriptions (fine-tuned); on the right, the model generating object-centric descriptions (pre-trained).

**Confidence scores**   We analysed the confidence score produced at generation time by the model for those captions which do not change after ablation, as this is an indicator of the extent to which the object ablation affected the generative process, even though it resulted in the same output sequence.

As shown in Figure 4.12, after ablation pre-trained VinVL generates object-centric descriptions with higher confidence than fine-tuned VinVL does with scene-level descriptions. However, the variance in the confidence score after ablation was lower for the fine-tuned model generating scene-level captions (Figure 4.13), suggesting greater robustness to ablation during scene-level caption generation.

**Bias check**   From the HL Dataset analysis in Section 3.4 we are aware that scene-types are not equally distributed in the HL-scenes data. As observed in Section 4.4.3, the sensibility to the ablation of diagnostic objects depends also on the scene-type represented in the image, thus we deem important to validate our results also on the train-val split rather than only on the test split, in order to avoid any potential distributional bias.

Therefore, we repeated the ablation experiment on both the test and the train-val

Figure 4.14: Kernel density estimate of the confidence scores distributions of unchanged captions after ablation for the test (blue) and train-val (orange) split.

split. The results obtained on the latter mirror those reported so far on the test-split only. In Figure 4.14 we show the comparison of the distributions of the unchanged confidence scores after ablation for the test and train-val split. There is no statistically significant difference between the distributions of confidence score shifts of the test set (shown in Figure 4.13) and the train-val set ($z = 0.13$ with $p = 0.89$ and $\alpha = 0.05$). This suggests that there are no significant distributional biases between train-val and test split such that could affect results in the scope of our analysis.

## 4.4.5 | Can we disentangle the role of attention and model representation?

The results so far suggest that there are significant changes in the model's self-attention when fine-tuning the model to scene-caption generation, though it keeps relying on diagnostic objects to generate scene-level captions. It is also somewhat more robust to object ablation, especially in the fine-tuned case. At this point, we probed the model's representations to address to what extent the knowledge required for scene-level caption generation was already present after pre-training. This would imply that the primary change to the model after fine-tuning is in the self-attention mechanism.

**Method**    Given a pair $(V, L)$ consisting of visual features $V$ and object labels $L$, we trained a probe to classify scene type based on VinVL encodings, before and after fine-tuning. We also repeated the procedure on inputs ablated as described in Section 4.4.4.

94

Figure 4.15: Scene distribution in the probing dataset

For this experiment, we identified 1426 images from HL-scenes, representing 8 types of scenes, down-sampling the most frequent classes, in order to have a ration between the least and most frequent class no less then 10%. The class distribution is shown in Figure 4.15. For every image in the probing dataset, we extracted the model's feature representations from the last layer and we averaged across the inputs, obtaining a single vector.

**Model selection**    We tested two probing models: a multi-layer perceptron and a random forest. We performed hyperparameter tuning of the neural probe by carrying out a random search followed by a probabilistic search. The tuned neural probe was a three-layer feed-forward network with hidden size 16, optimized using Limited-memory Broyden–Fletcher–Goldfarb–Shanno algorithm (LBFGS) with adaptive learning rate and $\alpha = 1$. Note that no parameter tuning was required for the random forest. As reported in Table 4.11, the random forest performs better or on par with the neural probe. Therefore we reported the performance of the random forest in the main results.

**Challenging the probe**    The probing model performed at ceiling with the more typical 90/10 split, especially when trained on the fine-tuned features (see Table 4.12). Therefore, we performed multiple experiments for different train/test splits namely, 90/10, 70/30, and 50/50. The 50/50 was the most challenging for the probe and it allows us to highlight the performance gap across different settings. Results from all the splits are shown in Table 4.12.

95

| Probe | Model | micro-F1 | macro-F1 | weighted-F1 |
|---|---|---|---|---|
| Random Baseline | | 0.16 | 0.12 | 0.16 |
| Random Forest | PRE | 0.94 | 0.67 | **0.92** |
| | FT | **0.99** | **0.96** | **0.99** |
| | PRE (A) | 0.92 | 0.66 | 0.90 |
| | FT (A) | 0.98 | **0.88** | 0.97 |
| Multilayer perception | PRE | 0.94 | 0.67 | 0.91 |
| | FT | 0.98 | 0.91 | 0.98 |
| | PRE (A) | 0.92 | 0.66 | 0.90 |
| | FT (A) | 0.98 | 0.85 | 0.97 |

Table 4.11: F1-scores of scene classification task in the 50/50 split. Models are trained on encodings extracted from the pre-trained (PRE) and fine-tuned (FT) model without and with ablation (A). In bold the best result for each setting.

| Split | Model | micro-F1 | macro-F1 | weighted-F1 |
|---|---|---|---|---|
| 90/10 | Pre-trained | 0.96 | 0.71 | 0.94 |
| | Fine-tuned | **1.0** | **1.0** | **1.0** |
| | Pre-trained (A) | 0.95 | 0.69 | 0.94 |
| | Fine-tuned (A) | 0.99 | 0.99 | 0.99 |
| 70/30 | Pre-trained | 0.94 | 0.67 | 0.92 |
| | Fine-tuned | **0.99** | **0.97** | **0.99** |
| | Pre-trained (A) | 0.93 | 0.66 | 0.91 |
| | Fine-tuned (A) | 0.98 | 0.94 | 0.98 |
| 50/50 | Random | 0.16 | 0.12 | 0.16 |
| | Pre-trained | 0.94 | 0.67 | 0.92 |
| | Fine-tuned | **0.99** | **0.96** | **0.99** |
| | Pre-trained (A) | 0.92 | 0.66 | 0.90 |
| | Fine-tuned (A) | 0.98 | 0.88 | 0.97 |

Table 4.12: F1-scores for scene classification task the random forest in different train/tes splits. The random forest is trained on encodings extracted from the Pre-trained (Pre-trained) and fine-tuned (Fine-tuned) model without and with ablation (A).

**Results**   We focus on the probe results for the 50/50 train/test split (see Figure 4.16), which is also the most challenging. The random baseline randomly assigns a label to the input features. For both pre-trained and fine-tuned models, probes perform at ceiling for scenes with a high support (cf. Figure 4.15). For scene types with a very low frequency, like *restaurant* and *room*, the probe trained on features from the pre-trained model fails. In contrast, probing features from the fine-tuned model still perform at ceiling. These results suggest that the information to detect the scene type is already present to some extent in the pre-trained model. Nevertheless, fine-tuning proves effective in closing the gap for low-support scenes.

When trained on features extracted from ablated inputs in Table 4.11, the probe is

Figure 4.16: F1-scores of the scene classification task for the pre-trained in (blue) and the fine-tuned model (orange).

not particularly affected by the ablation, confirming the robustness of the model's representations as observed in the ablation study (Section 4.4.4).

## 4.4.6 | Conclusions

In this study, we focused on scene-level caption generation. Taking a cue from prior work on scene semantics and syntax, our goal was to assess VL models' ability to reason about the link between scenes and their components and exploit this to generate informative captions with less redundancy.

Our analysis showed that the fine-tuning results in a more even distribution of attention mass over the image, suggesting a more "holistic" view of the scene which nevertheless makes use of diagnostic object information. Using a combination of ablation and probing methods, we also show that much of the relevant information for scene-level captioning is present after pre-training. Hence, the model's ability to generate scene-level captions is primarily acquired through a change in its self-attention.

# 4.5 | Summary

Motivated by research in cognitive science and human perception in scene understanding, in this Chapter we presented two studies focusing on analysing the capabilities of VL models in grounding scene-level descriptions.

In Section 4.3, we ran a first study where we tested these capabilities in VL encoders in zero-shot conditions. We found that these capabilities are not naturally learned by all the models. In fact, models seemed to be more sensitive to the style of the captions than to the actual content. However, a large-scale pre-training dataset and a constrastive learning objective seemed to play a role in the effective grounding of scenes. When scene-grounding is successful, the models exploit object-level visual clues to match the scene descriptions. This is performed by exploiting the visual information of diagnostic objects present in the scene. This finding is consistent with what has been observed in the human perception of scenes (e.g. Võ and Wolfe, 2013; Võ, 2021) and shows that stochastic models' optimization leads to object-level expectations in a visual scene that are similar to the humans ones.

After establishing that scene grounding capabilities are not naturally present in pre-trained VL models, we continue our investigation by studying the effect of the exposure of VL models to scene-level descriptions, in the generative setting. In the follow-up study presented in Section 4.4 we perform an in-depth analysis of the effect of standard fine-tuning of pre-trained VL captioning models, on scene caption generation. We found that fine-tuning on scene descriptions results primarily in a different allocation of attentional resources on the image, which consists in a more evenly distributed attention over the visual inputs. In other words, the model acquires a 'holistic' view of the scene without losing the capability to identify the single objects. In fact, the model keeps allocating more attention to the diagnostic objects relevant for the scene. This result is consistent with the findings reported in Section 4.3 were we link this aspect to the human perception of the scenes. Moreover, we do not observe any substantial change in the model's representations. In fact, the pre-trained representations show to be robust enough to support the majority of scene-types, albeit the successive fine-tuning helps strengthening the model's representations for weakly supported scene-types.

In this Chapter, we presented two studies focused on the exploration of the capabilities of VL models grounding scene descriptions. We exploit the alignment of object-level and scene-level captions, provided by the HL dataset (introduced in Chapter 3) to analyse scene grounding capabilities of VL models in terms of objects and their visual layout in the scene. Our method relies on well-known explainability techniques such as attention analysis, probing tasks and ablation studies. These methods allow a fine-

grained evaluation of specific aspects of the models however, they are model-specific, and require a specialized design and setup which are hardly scalable and generalizable to other scenarios. Model-agnostic XAI methods such as SHAP (Lundberg and Lee, 2017) are common in other ML domains, but difficult to apply to VL models, especially in generative settings, due to their high compute cost. Moreover, they are not designed to provide semantically informed explanation, exploiting abstract linguistic concepts, e.g. high-level captions. In Chapter 5 we will tackle these issues, proposing an explainability framework adaptable to these scenarios and general enough to be applied to any VL generative model, hoping to inspire future developments in this direction.

# How to explain High-level descriptions in VL generative models

The material in this Chapter is based on: Michele Cafagna, Lina M. Rojas-Barahona, Kees van Deemter, Albert Gatt. *Interpreting Vision and Language Generative Models with Semantic Visual Priors*, 2023. Frontiers in Artificial Intelligence Journal;

**Contributions**: Michele Cafagna: implementing and running the experiments, conducting the evaluation; writing and revising the paper. Albert Gatt and Lina M. Rojas-Barahona: supervising the research, writing, and revising the paper. Kees van Deemter: providing feedback with a particular focus on evaluation and revising the paper.

# 5.1 | Introduction

Multimodal learning research has witnessed a surge of effort leading to substantial improvements, in algorithms involving the integration of VL, for tasks such as image captioning (Hossain et al., 2019; Lin et al., 2014c; Sharma et al., 2020) and visual question answering (Antol et al., 2015; Srivastava et al., 2021; Zhu et al., 2016). The need has arisen to create more challenging tasks and benchmarks requiring higher fine-grained linguistic capabilities (Li et al., 2023a; Parcalabescu et al., 2022b; Thrush et al., 2022) and semantic and temporal understanding (Kesen et al., 2023; Park et al., 2020; Yu et al., 2016).

In this context, the role of interpretability methods has become central to assessing the models' grounding capabilities. In Chapter 4 we used several of these techniques, such as attention analysis, input ablation and probing tasks, to study and get some insight on the inner working of VL models when trained to generate scene descriptions, a kind of high-level captions introduced in the HL dataset (Cafagna et al., 2023b) described in Chapter 3.

However, such methods often need to be adapted for specific classes of tasks or models, lacking flexibility and generalization over new setups. To overcome this limitation, model-agnostic interpretability methods, such as SHAP-based methods (Lundberg and Lee, 2017), are often preferred, as they rely on a solid theory and benefit from desirable properties not available in other methods.



Figure 5.1: Example of token-by-token visual explanations using SHAP and superpixels features. A single visual explanation (heatmap) is generated for each generated token.

When such methods are applied to VL generative tasks, like image-captioning, the goal is to explain the textual output with reference to the visual input. However, the text generation process happens token-by-token, and as a result, most of the interpretability methods applied in this context tend to produce local token-specific explanations. Moreover, for most applications, current methods build the explanation on top of arbitrary regions of the visual input, usually considering superpixels (regions of adjacent pixels

of a fixed size) as the features against which to interpret the outputs (e.g. Parcalabescu and Frank, 2023).

Token-by-token explanations are hard to interpret as they are token-specific, and they are costly to compute since the number of model evaluations grows exponentially with the number of features used in each explanation. To mitigate these issues, approximation techniques, like sampling and input feature reduction, are usually applied. However, this produces inaccurate explanations which lack detail and are hard to interpret. An example is shown in Figure 5.1[1], with a visual explanation, namely a heatmap, highlighting the portion of the image affecting the token prediction computed to explain each single generated token using SHAP (Lundberg and Lee, 2017). While this kind of explanation is useful to explain single tokens referring to objects and entities in the image such as "bird" and "tree", it is unclear how to interpret explanations for tokens like "on" or "a" for which there is no direct connection with the image.

Furthermore, the reliance on superpixels as input features makes interpretation harder, since superpixels do not necessarily correspond to semantically meaningful regions of an image. In this Chapter, we address these issues by proposing:

1. A modular framework to create a new family of tools to generate explanations in VL generative settings;

2. A method to generate sentence-based explanations for vision-to-text generative tasks, as opposed to token-by-token explanations, showing that such explanations can efficiently be generated with SHAP by exploiting semantic knowledge from the two modalities;

3. A method to reduce the number of visual input features by exploiting the semantics embedded in the models' visual backbone. We extend this method to a number of different architectures. We further propose an alternative approach to extract semantically meaningful features from images in case a model architecture does not support our specific method;

4. A human evaluation designed to assess key user-centric properties of our explanations.

---

[1]Image from shap.readthedocs.io

# 5.2 | Related Work

In this section, we discuss related work on interpretable machine learning and XAI in Vision and Language models. We also detail some of the essential properties of the XAI framework (SHAP) which our work is based on.

## 5.2.1 | Interpretable Machine Learning

Interpretable machine learning is a multidisciplinary field encompassing efforts from computer science, human-computer interaction, and social science, aiming to design user-oriented and human-friendly explanations for machine learning models. It plays an important role in the field for a series of reasons: it increases trust, confidence, and acceptance of machine learning models by users, and enables verification, validation, and debugging of machine learning models. As discussed in Section 2.3.4 explainability techniques for DNN can be grouped into two main categories: *white-box* methods which exploit the knowledge of the internal structure of the model to generate the explanation and *black-box* methods, also called model-agnostic, which operate only on the inputs and the outputs (Loyola-Gonzalez, 2019).

**White-box methods**    There exist two types of white-box methods: attention-based and gradient-based methods.

   *Attention-based* methods (e.g Ahmed et al., 2021; Zheng et al., 2022) exploit the model's attention activations to identify the part of the input attended by the model during the prediction. As shown, in Section 4.4.3, they can be used to explain predictions in the image captioning as well as other tasks like image recognition (Li et al., 2021c), authorship verification (Boenninghoff et al., 2019) gender bias identification (Boenninghoff et al., 2019) etc.

   On the other hand, *Gradient-based methods* (e.g. Selvaraju et al., 2017; Springenberg et al., 2014) compute feature attributions by manipulating the gradients computed in the backward step with respect to the original inputs (Shrikumar et al., 2016), or with respect to a specific baseline (Simonyan et al., 2013; Sundararajan et al., 2017).

**Black-box methods**    do not make any assumptions regarding the underlying model. For example, Permutation Feature Importance (Breiman, 2001), initially designed for random forests and later extended into a model-agnostic version by Fisher et al. (2019), consists in randomly shuffling the input features and evaluating the model's output variations. Ribeiro et al. (2016a) proposed LIME (Local Interpretable Model-Agnostic

Explanation), which uses a surrogate linear model to approximate the black-box model locally, that is, in the neighborhood of any prediction. LOCO (Lei et al., 2018) is another popular technique for generating local explanations. It can provide insight into the importance of individual variables in explaining a specific prediction. SHAP (Lundberg and Lee, 2017) is a framework considered by many to be the gold standard for local explanations, thanks to its solid theoretical foundation. SHAP leverages the concept of Shapley values, first introduced by (Shapley et al., 1953), used to measure the contribution of players in a cooperative game. This was later extended by (Lundberg and Lee, 2017) for the purpose of explaining a machine learning model.

In this Chapter, we propose a flexible hybrid framework based on SHAP, which benefits from properties typical of *black-box* methods, since it can be applied in a completely model-agnostic way. At the same time, our method shares some properties with *white-box* approaches since, when possible, it takes advantage of certain internal components of the model. In particular, the framework we propose for VL generative models can be leveraged to exploit architectural features of a model's visual backbone to generate more semantically meaningful explanations.

## 5.2.2 | Background on SHAP

In the context of machine learning, the cooperative framework introduced by Shapley et al. (1953) can be framed as a game where each input feature is a player and the outcome is determined by the model's prediction. Shapley values measure the contribution of each player to the final outcome, or in other words, the input features' importance. Shapley redistributed the total outcome value among all the features, based on their marginal contribution across the possible coalitions of players, i.e. combinations of input features. The outcome of the game, namely the prediction of the model, is redistributed across the features, in the form of contributions that have three desirable properties:

- *Efficiency*: all the Shapley values add up to the final outcome of the game;

- *Symmetry*: all the features generating the same outcome in the game have the same Shapley value, thus the same contribution;

- *Dummy*: if adding a feature to a coalition (i.e. set of features) does not change the outcome of the game, its Shapley value is zero.

The Shapely values compute cost grows exponentially with the number of players, i.e input features, used in the game, as it requires the model's evaluation of all the possi-

ble combinations of features. To tackle this issue, Lundberg and Lee (2017) contribute by formulating a variety of methods to efficiently approximate Shapley values in different conditions:

1. KernelSHAP: derived from LIME and totally model agnostic, hence the slowest within the framework;

2. LinearSHAP: designed specifically for Linear models;

3. DeepSHAP: adapted from DeepLift (Shrikumar et al., 2017) for neural networks, which is faster than KernelSHAP, but makes assumptions about the model's compositional nature.

Later on, the framework was extended with other methods with variations for specific settings; Mosca et al. (2022b) propose a thorough description of the SHAP family of methods.

It is important to note that all these methods work under the so-called *feature independence assumption*, which is fundamental for the theoretical resolution of the problem. The feature independence assumption says that the features in the SHAP game do not correlate or overlap with each other. Since Shapley (and SHAP) attributions are computed by marginalising features, if a feature is strongly correlated to or overlaps with another one, such marginalisation yields unrealistic results (Molnar, 2020). However, in order to deal with real-life scenarios, this constraint can be relaxed to some extent. For instance, in NLP tasks each token of a textual sequence is considered an independent feature (Kokalj et al., 2021) whereas, in Computer Vision, the image is usually split into squared patches or superpixels, which are also considered independent of each other (Jeyakumar et al., 2020). In both of these cases, the independence assumption is a simplification. For example, language tokens are often mutually dependent in context (and this is indeed the property leveraged by self-attention in Transformer language models). Similarly, pixels in neighboring patches in an image may well belong to the same semantically relevant region (and this is indeed the property exploited by neural architectures suited for computer vision tasks, such as convolutional networks or vision transformers). Properties of tokens in context and those of pixels in image regions have been taken into account in some adaptations of SHAP which consider the hierarchical structure of the feature space, such as HEDGE for text (Chen et al., 2020a) and h-SHAP for images (Teneggi et al., 2022).

Along the same line, in our work, we relax the independence assumption providing in Section 5.4.4.2 a detailed analysis and discussion of this issue.

## 5.2.3 | Kernel Shap

The core method of our framework is Kernel Shap. We base our approach on the formulation by Lundberg and Lee (2017), which provides an accurate regression-based, model-agnostic estimation of Shapley values. The computation is performed by estimating the parameters of an explanation model $g(x')$ which matches the original model $f(x)$, namely:

$$f(x) = g(x') = \phi_0 + \sum_{i=1}^{M} \phi_i x'_i \tag{5.1}$$

where $M$ is the number of input features (or players) and $x'_i$ is a player of the game. $g(x')$ is approximated by performing a weighted linear regression using the Shapley kernel:

$$\pi_{x'}(z') = \frac{M-1}{(M \text{ choose } |z'|)|z'|(M-|z'|)} \tag{5.2}$$

where $z'$ is the subset of non-zero entries, namely a binary representation of the coalition of players. The Shapley kernel, in other words, is a function assigning a weight to each coalition. The number of coalitions needed to approximate the Shapley values corresponds to all the possible combinations of players, i.e. $2^M$ coalitions. This makes Kernel SHAP extremely expensive to compute (and slow in practice) when $M$ is large.

Our framework relies on KernelSHAP as is it totally model-agnostic. We address both the efficiency issue and the strict independence assumption of the method by generating semantic input features (more details in Section 5.3.2.2) and optimizing the approximation through sampling (full details in Section 5.3.1).

## 5.2.4 | Explainability for Vision and Language

One way to characterize the scope of VL models is with respect to the types of tasks they are designed to address. On the one hand, tasks like image captioning (Anderson et al., 2018; Fisch et al., 2020b; Li et al., 2022c; Mokady et al., 2021; Zhang et al., 2021), image-text retrieval (Cao et al., 2022; Radford et al., 2021b), and visual question answering (Antol et al., 2015) require a strong focus on the recognition of objects in images. More recently, research has begun to explore the capabilities of models in tasks that require some further reasoning or inference over the image contexts, such as understanding analogies (Zhang et al., 2019a), describing actions and rationales (Cafagna et al., 2023b) and inferring temporal relations (Park et al., 2020).

The need to understand how VL models ground their predictions has become essential, leading to the emergence of Explainable Artificial Intelligence (XAI) for multimodal

settings (Zellers et al., 2019). Visual explanations can help humans to know what triggered the system's output and how the system attended to the image. To this purpose, feature attribution methods are often preferred as they can provide a visual explanation of the prediction. Most of the XAI methods introduced for unimodal tasks can be adapted to VL tasks.

Some popular *white-box* methods use gradients to generate saliency maps to highlight the pixels corresponding to highly contributing regions. These methods include Grad-CAM (Selvaraju et al., 2017; Shrikumar et al., 2016) or Layer-wise Relevance Propagation (LRP) (Binder et al., 2016) where the contribution is computed with respect to an intermediate layer instead of the input layer. These methods can produce fine-grained pixel-level explanations. However, their outcomes can be noisy and require many evaluations to converge to a stable explanation.

*Black-box* approaches are mostly perturbation-based, that is, they compute attributions based on the difference observed in the model's prediction by altering the input. Such methods include occlusion sensitivity (Uchiyama et al., 2023), RISE (Petsiuk et al., 2018), and LIME (Ribeiro et al., 2016a). Other approaches are task-agnostic, like MM-SHAP (Parcalabescu and Frank, 2022), where a SHAP-based method is used to measure the contribution of the two modalities in VL models independently of the task performance. Although these methods make few assumptions about the underlying model, their explanations are computationally expensive, as the number of model evaluations required grows exponentially with the number of features. To overcome this limitation, the number of features is usually reduced by partitioning the image into patches called superpixels, which discretize the input into a smaller number of features. However, this approach can lead to coarser and not very informative explanations.

Explanations for VL **generative tasks**, like image captioning, incur even more complexity, as the prediction of the model is now a textual sequence. As noted in Section 5.2.2, SHAP estimates feature contributions based on the amount of variation observed in the model output, with or without the feature. This requires a numerical output value (which of course, linguistic sequences are not). A popular solution, which is in keeping with the autoregressive nature of neural language decoders, is to break down the caption generation process into a series of steps where each token is explained separately with respect to the image and the previously generated sequence. This requires generating a single visual explanation for each generation step. However, the meaning of the sentence is not only determined by the meaning of the single words it is composed of but also by the way these words are combined and arranged together. Therefore, a global meaningful explanation must take into account the whole textual sequence and not just part of it, as only in this way can the explanation take into account the whole

textual context.

A popular solution is to generate the token-level explanations using Integrated Gradients (Sundararajan et al., 2017), providing region-level visualizations or using the attention activation scores to visualize the model's attended regions (Cornia et al., 2022; Zhang et al., 2019b). However, these methods are white-box approaches as they make assumptions about the inner workings of the model; thus they need to be specifically re-adapted to new systems. Furthermore, they focus on token-level explanations but do not allow a comprehensive global explanation of the textual output.

To the best of our knowledge, our work is the first attempt to bring together a model-agnostic framework like SHAP, in the image-to-text task, with the aim of providing a comprehensive explanation of the generated textual output as a whole, rather than on a token-by-token level.

We further propose a method to provide explanations based on features that are semantically meaningful, rather than on patches or superpixels.

## 5.3 | Method

In this section, we first address the matter of efficiency which, as noted above, is a pressing problem for methods based on Kernel SHAP. We then turn to the core proposals in our method, adapting it to generative models to achieve explanations for whole sequences rather than tokens (Section 5.3.2.1) and using semantically meaningful visual regions as features (Section 5.3.2.2).

### 5.3.1 | Deterministic Kernel SHAP sampling

Kernel SHAP is model-agnostic, meaning that it cannot make any assumption on the model to explain. For this reason, it is also among the slowest in the SHAP family of XAI methods (Mosca et al., 2022a). This issue is addressed by performing Monte Carlo sampling over the pool of coalitions, allowing under certain conditions to compute a reasonably accurate approximation of Shapley values, even in the case of large-sized models or low-resource hardware.

Taking inspiration from Molnar (2020), we implemented a deterministic sampling strategy. Given a specific sampling budget $k$, we prioritized coalitions which have a high weight, where weight is computed by Eq. 5.2. This was achieved by generating the coalitions in decreasing weight order and selecting the first $k$ coalitions. In Figure 5.2, we compare the weights of coalitions computed using the standard Kernel SHAP (on the left) and using the method which prioritises high-weight coalitions (on the right).

Figure 5.2: Standard Kernel SHAP (left) and modified Kernel SHAP with priority for high-weight coalitions (right). The y-axis corresponds to weight whereas the x-axis is the iteration in which a particular coalition is generated.

As can be observed, our sampling strategy with priority (on the right) ensures that we select the high-weight coalitions first, providing an optimal ordering among samples.

Our experiments show that sampling with priority offers two main advantages:

1. higher accuracy of the Shapley values estimate;

2. a deterministic sampling strategy.

In Figure 5.3 we report the approximation error of the Shapley values when applying Kernel SHAP, using Monte Carlo (orange) and the deterministic high-weight priority (blue) as sampling strategies, for different sample sizes. The error is computed over 10 runs, using the Mean Squared Error (MSE) with respect to the Shapley values computed with Kernel SHAP using all the $2^M$ coalitions. Our deterministic sampling approximates Shapley values with errors that are orders of magnitude smaller than Monte Carlo sampling. We observe this consistently for different sampling sizes.

With a more efficient and deterministic sampling strategy, we now turn to the core of our method.

## 5.3.2 | Adapting Kernel SHAP to vision and language generative tasks

In the image captioning scenario, we can set up a cooperative game, where we want to compute the contributions of the players, i.e. the pixels of the image, with respect to the outcome, i.e. the caption. In Section 5.2, we identified two shortcomings of the standard way in which this is performed. Here, we discuss our contributions to overcome these shortcomings.

The first problem is related to the comprehensiveness of explanations. In order to measure the variations of the outcome of the function needed to run Kernel SHAP, the

Figure 5.3: Mean Squared Error (MSE) of the Shapley values estimated using Monte Carlo sampling (orange) and deterministically sampling coalitions with high priority (blue), for various sampling sizes. All the values on the x-axis are exponentials ($2^{M-1}, 2^{M-2}, 2^{M-3}$) where $M$ corresponds to the number of features. The MSE is computed with respect to the Shapley values computed using all the $2^M$ coalitions available in the sampling space.



Figure 5.4: Overview of the explainability framework. The new components proposed in this work are shown with a dark border. Our method leverages **KernelSHAP** as the core explainability method. We introduce semantic features extracted using **DFF** from the captioner's visual backbone and generate **sentence-based** visual explanations based on the estimated Shapley values.

111

caption generation process is usually broken down into token generation steps. Each step produces logits that can be used to compute a numerical outcome. However, this forces us to consider each generation step as a separate cooperative game, meaning that we need to run a separate instance of Kernel SHAP for each generated token, further scaling the time and compute cost needed to explain an image-caption pair, linearly in the length of the generated sequence. Moreover, such explanations refer to single tokens and do not provide an explanation for the whole output of the model, namely the caption.

The second problem is related to the definition of coalitions in the visual input. The number of coalitions to be computed grows exponentially with the number of players. This makes the computation of the Shapley values intractable for images, as their basic features are pixels. Therefore, the huge number of coalitions makes any sampling strategy inaccurate when considered reasonable sampling budgets. In order to overcome this limitation, the image is typically partitioned into a grid composed of *superpixels*, namely groups of pixels, each of which represents a single player. This reduces the total number of players in the game, making computation of the Shapley values more feasible, but at the same time, it reduces the degree of the detail of the explanation. Moreover, we argue that dividing the image into a grid of square superpixels breaks the semantics underlying the image, resulting in potentially under-informative explanations. In particular, there is no guarantee that the pixels grouped together in this manner correspond to semantically meaningful image regions.

In the following sections, we address these issues, proposing alternative solutions. Specifically, we address the first shortcoming in Section 5.3.2.1, before turning to a proposal for semantically meaningful and sparse features in Section 5.3.2.2. Our solution can be integrated with existing methods, to compose a modular explainability framework for generative VL models. An overview of this framework is shown in Figure 5.4.

### 5.3.2.1 | Towards sentence-based explanations

In order to adapt Kernel SHAP to generate global explanations for the caption, we measured variations of the caption's meaning representation when perturbations are applied to the input image. This allowed us to numerically quantify the meaning variation of the whole caption which is due to the marginal contributions of different input features (image regions or pixels).

Formally, given an image-captioning model $f$ and an image $x$ we generated a caption $c = f(x)$ and we computed:

$$e_{ref} = E(c) \tag{5.3}$$

Figure 5.5: Example of the sentence-based explanation. 1) We compute the reference embedding (red) from the caption generated by the model when the input has no perturbation. For each perturbation applied, we compute the embedding (orange, blue) of the resulting caption and use the cosine distance between the reference and the current embedding, to measure the semantic variation of the caption.

where $e_{ref}$ is the embedding representation of $c$ that we considered the *reference embedding* of the caption, and $E()$ is a function used to extract such a representation.

For each perturbed image $x'$ and its corresponding caption we extracted, analogously, an embedding $e'$. Then we computed:

$$s = cos(e_{ref}, e') \qquad (5.4)$$

where $s$ is the variation in the embedding representation computed as the cosine distance $cos(\cdot)$, between the reference embedding $e_{ref}$ and the embedding of the caption of the perturbed image $e'$.

In other words, we used the cosine distance between the semantic representation of the reference caption and the caption generated upon input perturbation, to measure the model's output variations. A schematic representation of the method is shown in Figure 5.5.

Re-framing the problem as described, allowed us to apply Kernel SHAP to compute feature attributions taking into account the semantic variation, i.e. the cosine similarity between the original and the perturbed caption, of the whole caption in a single cooperative game instance.

### 5.3.2.2 | Exploiting semantic visual priors

Partitioning the image into a grid of superpixels is a straightforward way to reduce the number of input features in the image. As argued above, although convenient, superpixels do not guarantee the preservation of semantic information depicting the visual content, as they shatter the image into equally sized patches regardless of the content represented. We addressed this issue by proposing a semantically guided approach, that selects the input features according to semantics-preserving visual concepts arising from the visual backbone of the VL model.

This not only allows for generating more meaningful explanations but explicitly focuses explanations of the model's generative choices on the output of the model's own visual backbone.

We generated input features leveraging the Deep Feature Factorization (DFF) method (Collins et al., 2018). DFF is an unsupervised method allowing concept discovery from the feature space of CNN-based models. We refer to such concepts as 'semantic priors', that is, the knowledge or assumptions learned by the visual backbone, in the context of a given domain or task. We used them to craft input features that produce semantically informed visual explanations.

Formally, following Collins et al. (2018)'s notation, given the activation tensor for an image $I$: $A \in \mathbb{R}^{h \times w \times c}$ where $h, w, c$ correspond respectively to the height and width, and the number of channels of the visual backbone's last activation layer, we performed a Non-negative Matrix Factorization (NMF) of $A$:

$$NMF(A,k) = \underset{\hat{A}_{Ik}}{\arg\min} \|A - \hat{A}_k\|_F^2,$$

$$\text{subject to} \quad \hat{A}_K = HW, \forall i, j : H_{ij}, W_{ij} \geq 0, \tag{5.5}$$

where $W \in \mathbb{R}^{n \times k}$ and $H \in \mathbb{R}^{k \times m}$ enforce the dimensionality reduction to rank $k$.

Each column $H_j$ was reshaped into $k$ heatmaps of dimensions $h \times w$, each of which highlights a region which the factor $W_j$ corresponds to. The heatmaps were then up-sampled to match the original image size with bilinear interpolation and converted into binary masks, each of which corresponding to an input feature. In this way we obtained $k$ input features, where $k$ is the number of regions extracted. A schematic example of input feature extraction performed by DFF is shown in Figure 5.6.

In our method, the regions identified via DFF are the features for which attributions are computed. The key intuition is that these features correspond to meaningful sub-parts of the input image according to the VL model's visual backbone. They do not necessarily reflect humans' visual expectations of the image (although we find that

Figure 5.6: Schematic example of input features extraction using DFF. Through thresholding, we convert the heatmaps into binary masks that we use to create semantically meaningful features.

they often do); rather they represent the visual priors learned by the vision model after training.

To create a coalition we summed up multiple masks, then applied them to the original image, which later contained only pixels belonging to input features in the selected coalition.

NMF can be seen as an unsupervised clustering algorithm, allowing control for the number of clusters or concepts to find. $k$ can be considered a hyperparameter of the method, which we show can be kept small to achieve a good level of semantic detail and low compute cost.

**Non-partitioning features**   DFF generates semantic masks reflecting the activations of the model's visual backbone. The whole process is unsupervised and produces masks that do not constitute partitions of the image, meaning that it is not guaranteed that the sum of all the extracted masks will match the total size of the image.

In order to account for this issue, we created an additional *leftover* mask covering the remaining area and we included it in the SHAP cooperative game. This allowed us to consider the whole visual information represented by the image, in the game. As noted in Section 5.2.2, the computation of Shapley values is based on a feature independence assumption. Since our features may be non-partitioning, this constraint may not hold, thus we relaxed this assumption in our approach. We explored the consequences of this in more detail in Section 5.4.4.2.

**Intensity-preserving explanations**   SHAP-based methods relying on superpixels assume that each pixel in a patch contributes equally, thus all the pixels in a patch are assigned the same Shapley value. However, in DFF, features in each binary mask correspond to an equally sized heatmap. Therefore, we multiplied the Shapley value by the heatmap corresponding to the binary mask. This allowed us to exploit the models'

visual priors by scaling the contribution according to the intensity of the feature signal. In other words, we used the intensity of the visual backbone's activations to highlight regions of the image within the input feature according to the Shapley value estimated for that input feature. This resulted in a more fine-grained and visually detailed explanation, without additional compute cost.

# 5.4 | Experiments

The methodology described in the previous section raises an important question which we now address experimentally: *What are the pros and cons of our method based on visual semantic priors in comparison with standard feature selection methods used in VL, based on superpixels?*

In this section, we describe the data and task, as well as a SOTA vision-to-language model, which we used to perform a human evaluation of our explainability framework.

## 5.4.1 | Data

We validated the method presented in the previous section with experiments using the HL dataset (Cafagna et al., 2023b) introduced in Chapter 3.

The systematic alignment in the HL dataset of the object-centric and abstract captions along three axes, i.e. *scene*, *action* and *rationale*, provides us with a suitable test bed to compare the efficacy of our method in delivering global explanations in both captioning and visual question-answering scenarios. As already discussed in Chapter 3, differently from the object-centric ones, high-level captions do not explicitly mention objects visually present in the scene. In fact, they use more abstract terminology, e.g. *They are having fun*, which often requires the evaluation of the whole sentence to obtain a meaningful explanation. An example pairing the three high-level captions and the

| Image | Axis | Caption |
|---|---|---|
|  | scene | at a sport field |
| | action | they are playing a sport |
| | rationale | they are having fun |
| | object-centric (COCO) | A woman has fallen on the ground in a field. |

Table 5.1: Example of High-Level captions. It is shown one of the three captions available for the three axes collected: *scene, action, rationale*, aligned with the object-centric captions from COCO.

original COCO caption from the HL Dataset is shown in Table 5.1. For full details see Section 3.4.

One of the distinguishing features of our method is the exploitation of the model's visual priors to generate visual explanation. Therefore, we are also interested in experimenting on the capability of our method to provide visual explanations for different kinds of description. This makes, the visual question answering task an interesting scenario to test our method on, as it allows us to observe and compare visual explanations computed on different kinds of description, using the very same model, namely without optimizing a single model for each kind of description, as instead done in Section 3.5. Moreover, this setup is perfectly suitable to the HL dataset, given that the high-level captions have been collected using a question-answering setting.

## 5.4.2 | Model

For our experiments, we focused on one VL model, since our goal was to evaluate the quality of explanations, not the model itself. Our choice was motivated by two considerations: first, a model should ideally have good performance in zero-shot settings; second, it should exhibit SOTA performance on generative tasks. OFA (Wang et al., 2022b) is a large pre-trained multimodal model with a CNN-based visual backbone, trained using a task-agnostic and modality-agnostic framework. OFA is able to perform a diverse set of cross-modal and unimodal tasks, like image captioning, visual question answering, image generation, image classification, etc. It is trained on a relatively small amount of data (20M image-text pairs) with instruction-based learning and a simple sequence-to-sequence architecture. Nevertheless, on downstream tasks, it outperforms or is on par with larger models trained on a larger amount of data. OFA is effectively able to transfer to unseen tasks and domains in zero-shot settings, proving to be well grounded also in out-of-domain tasks.

This makes OFA an excellent candidate to test our explainability framework in a real-world scenario, namely a large pre-trained generative model with state-of-the-art performance on downstream tasks in zero-shot conditions. Thus, we used OFA to generate textual predictions in a VQA setting. We then used our framework, which combines DFF features and sentence-based explanations, to generate visual explanations of such predictions. In our evaluation, we compared these explanations to the more standard setup to vision-to-text models, that is the one based on superpixels as features.

Figure 5.7: Global visual explanation for the question "What is the subject doing?", and corresponding model's answer "drinking". Explanations are generated using Kernel SHAP. The explanation using DFF input features (on the left) provides a detailed positive (blue) area. We use 11 DFF features and 12 superpixel features. The explanation generated by superpixel input features (on the right) although covering a similar region, i.e. the glass, does not provide the same level of detail.

## 5.4.3 | DFF vs Superpixel

In this Section, we focus on the comparison between the global visual explanations produced using superpixel or DFF input features. We focus on the capability of the two methods to adapt to different semantic aspects of the explanation; in Section5.4.3.1 we specifically address this discussion with a focus on the VQA task.

All the experiments were performed in zero-shot by using the *OFA-large* model in its original implementation [2]. In order to ensure a fair comparison, we extracted a similar number of features for both methods, namely 12 for superpixel and 11 for DFF. This number allowed us to execute the experiments in a reasonable amount of time. In fact, we recall that the number of features has an exponential impact on the number of model evaluations needed to generate the explanations. Reducing the number of features mitigates the efficiency issue, but does not solve it. An in-depth discussion about the efficiency issue is provided in Section 5.4.3.2.

As an initial comparison, Figure 5.7 shows a direct comparison between the two kinds of input features for the caption "drinking", generated using sentence-based Kernel SHAP. Both methods assign a positive contribution to the region corresponding to the glass, with some important differences:

■ **Detail**: The DFF features succeed in capturing the key visual semantics of the image, i.e. the glass, in a single input feature (with some noise), producing a more

---

[2]`https://github.com/OFA-Sys/OFA`

118

Q: *Where is the picture taken?*
A: *in a living room*

Q: *What is the subject doing?*
A: *eating*

Q: *Why is the subject doing it?*
A: *The subject is going to eat the pizza*

Figure 5.8: Examples of explanations for the VQA task from the HL Dataset for the *scene* (far left), *action* (center) and *rationale* (far right) axes. The top row shows the questions (Q) and the generated answers (A). The middle and the bottom row, show visual explanation generated respectively with DFF and superpixel input features, with comparable compute cost.

detailed explanation than superpixels, where the region corresponding to the glass is shared across different patches (i.e. different features).

- **Intensity**: DFF scales the contributions according to the magnitude of the feature signal (as described in Section 5.3.2.2), providing a fine-grained visual indication of the importance of specific sub-regions within the same input feature region.

### 5.4.3.1 | Semantic visual features improve the quality of the explanations

We compared DFF and superpixel explanations on the VQA task. We selected images and questions for the three axes in the HL dataset, i.e. actions, scenes, and rationales, and we generated visual explanations for the answers. This allows us to compare how the two methods handle semantically different aspects highlighted in the visual content.

We expected to see that the positive contribution assignment (in blue) changes for the same image for different captions, corresponding to different kinds of questions for

|            |            |              |
| :--------: | :--------: | :----------: |
| (a) 4X4    | (b) 8X8    | (c) 16X16    |

Figure 5.9: Example of explanations generated with superpixels, with an increasing number of features, namely $16, 64, 256$ features (respectively 5.9a, 5.9b, 5.9c). These are obtained with Kernel Shap sampling using a fixed sampling budget of 2048 samples.

which the model generates different answers. In response to different questions about location, rationale, or action, the model's output should depend on different regions of the image. For instance, we expected to observe a wider positive area highlighted in the picture for the *where* question and a more specific detailed area for the *what* question. As shown in Figure 5.8, the DFF-based method (first row) succeeds in highlighting in significant detail the semantic areas contributing to the output. On the other hand, superpixels (second row) provide coarser detail, as they are limited by the size of the patches. This suggests that the DFF-generated explanations could lead to a visible advantage in terms of comprehensiveness and completeness; we further test these hypotheses by running a human evaluation, in Section 5.5.

### 5.4.3.2 | Semantics-guided explanations are efficient

In order for superpixel-based explanations to achieve a level of detail comparable to DFF, we need to significantly increase the number of patches. However, this causes an exponential surge in computing cost, which makes it unfeasible to run, especially if we are testing large models. This issue can be mitigated by performing deterministic Kernel SHAP sampling (as described in Section 5.3.1). Combining the exponential growth of the sample space, and the limited sampling budget can easily lead to unreliable explanations. An example is shown in Figure 5.9 where we perform Kernel SHAP sampling using superpixel's features at increasing number of patches (i.e. $16, 64, 256$) in order to increase the detail of the explanation (i.e. smaller patches). In order to make a fair cost comparison we keep a fixed sampling budget of 2048 samples, which is the same budget used to compute the DFF explanation in Figure 5.7.

Positively contributing regions, corresponding roughly to the glass in Figure 5.9a,

Figure 5.10: Binary feature masks extracted using DFF with $k = 10$. The $11^{th}$ feature is the *leftover* mask. The original image is the same shown in Figure 5.7 and Figure 5.9.

change inconsistently for Figure 5.9b and 5.9c, due to the exponential growth of the feature combination space, resulting in unreliable explanations.

On the other hand, DFF does not suffer from this issue. In fact, there is no clear advantage in increasing the number of features, because the main semantic content is usually embedded in a small number of features. In our experiments, we established that a good number of features for DFF is between 8 and 12. This number of features keeps the computational cost low, allowing us to compute full Kernel SHAP or Kernel SHAP sampling with very high accuracy.

## 5.4.4 | Semantic features analysis

The semantic features extracted by DFF are drastically different from superpixel features in many key aspects related to the visual content captured. Moreover, DFF is unsupervised and dynamically exploits the visual backbone's priors. In this Section, we focus on analyzing the benefits and limitations characterizing the semantic features generated by DFF. We discuss in detail key aspects like the kind of semantic content captured along with possible theoretical implications and how it can be generalized over different visual backbones.

(a) Overlapping features          (b) Non-disjoint features          (c) Disjoint features

Figure 5.11: Example of overlap (highlighted in red) between two feature masks (Figure 5.11a) and comparison between visual explanations generated given the question "What is the subject doing?" and the model's answer "drinking". We compare regular DFF features (Figure 5.11b) and disjoint DFF features (Figure 5.11c). Although the masks overlap only to a small extent, the explanation is visibly affected.

### 5.4.4.1 | What kind of semantics do DFF features capture?

DFF features capture semantic concepts learned by the model's visual backbone. These do not necessarily follow human visual expectations. In Figure 5.10 we show an example: features 1, 2, and 8 can be associated with three main **semantic objects and entities** of the image shown in Figure 5.7, namely *face, glass* and *shirt*. However, we observe in the remaining features **several geometrical patterns**, that highlight the edges and the corners of the pictures. This pattern is recurrent in the features extracted by DFF, independently of the visual content. We believe this is partially due to the capability of CNNs to capture spatial configuration (Zeiler and Fergus, 2014) and the effectiveness of DFF in factorizing together model activations with similar characteristics.

### 5.4.4.2 | Relaxing the feature independence assumption

As described in Section 5.2.2, SHAP in the cooperative game formulation assumes the *feature independence principle*, namely that each feature is independent of all the others. However, this assumption does not hold for image data since each pixel is inherently dependent on the other pixels, especially those in its vicinity. Therefore, in order to work with visual data, this constraint needs to be relaxed. This solution is typically applied for computer vision tasks by graphical models like Conditional Random Fields (CRF). CRFs relax the strong independence assumption on the observations (the pixels of the image) by modeling the joint distribution of observations, usually intractable, as a conditional distribution (Li et al., 2022d).

   Along the same lines, superpixel features relax this constraint by partitioning the

image into patches that are not independent, considering the underlying semantics depicted in the visual content.

This issue is mitigated by the DFF features, as they tend to cover semantically related regions of the image, preserving the underlying visual semantics. On the other hand, as pointed out in Section 5.3.2.2, DFF features are not disjoint, meaning that to some extent, the contribution of overlapping regions is subject to contamination from other regions. In this section, we analyse the consequences of this in more detail. Our analysis follows two steps:

1. We measured the DFF feature overlap over a sample of 1000 images; finding that the amount of overlap among the feature masks corresponds to 0.77% of the pixels in the image with a standard deviation of 0.63 and an average maximum peak of 2.04%. This suggests that this phenomenon is present to a limited extent, at least for the model we are using.

2. We compared visual explanations generated by disjoint and non-disjoint features. In order to generate disjoint features, we post-process the feature masks extracted, by checking all possible pairs of feature masks and assigning the possible overlapping region to one of the two compared features. An example is shown in Figure 5.11a where the overlapping regions (highlighted in red) between two feature masks are randomly assigned to one of the features (either blue or green).

Enforcing the features' disjointness leads to similar results to their non-disjoint counterpart. However, in some cases, the re-allocation of the overlapped region impacts the Shapley value of the feature, causing unpredictable results. This suggests that manually changing the feature masks can disruptively affect the visual semantics captured by the feature, leading to misleading visual explanations. A cherry-picked example is shown in Figure 5.11, where using the disjoint features (Figure 5.11c) causes a substantive change in the visual explanation.

In conclusion, we observe that **the phenomenon of non-disjoint features is present to a limited extent** and overall **it does not invalidate the visual explanations**, as it can be considered a relaxation of the feature independence assumption. Moreover, as empirically observed, relaxing this assumption is unlikely to invalidate the method, as the explanation is consistent with the ones generated by superpixel features. On the other hand, we observe that **forcing the feature masks' disjointness harms their capability to preserve the visual semantics, leading to misleading visual explanations**.

Figure 5.12: RBO scores computed between normalized and unnormalized Shapley values, for positive (blue), negative (orange), and all (green) features.

### 5.4.4.3 | Does feature size matter?

Different from superpixel patches, DFF semantic features can have different sizes, depending on the semantic role of the highlighted region. We ask to what extent the size of a visual feature could affect the final contribution in the SHAP cooperative game. In order to test for that, we normalized the Shapley value obtained according to the size of the feature mask and we compared normalised values with the unnormalized ones. To normalize a Shapley value we computed:

$$r_i = \frac{\sum_{j=0}^{|M_i|} m_j}{|M_i|}$$
$$\hat{a}_i = \frac{a_i}{r_i}$$

(5.6)

where $m_j$ is a non-zero element of the binary mask, $|M_i|$ is the total number of entries in mask $i$ and $r_i$ indicates the proportion of the image covered by the mask. $r_i$ is then used to discount the magnitude of the Shapley value $a_i$ obtaining the normalized value $\hat{a}_i$.

In the normalization process, the feature contribution's magnitude is obviously rescaled. However, we are interested in measuring to what extent the normalization has affected the features' importance in relative terms. Therefore, we used the RBO (Webber et al., 2010), a similarity metric for ranked lists, to measure the difference in the feature attribution ranking after normalization for a sample of 100 DFF-based explanations. A significant change in feature ranking would entail a positive correlation between size and feature importance.

In Figure 5.12 we show the results of this experiment: the RBO is overall at ceiling, with a minimum value, including outliers, greater than 0.9 (in a range where 1 is identical ranking and 0 is totally different). The positive contributions, which are the most informative to understand the explanations, are the most stable in terms of ranking. This suggests that **the size of the features extracted using DFF does not significantly**

Figure 5.13: Schematic example of how to generate semantic features with DFF from a ViT visual backbone. The index of the highlighted band in the heatmap is used to select the patches to create the feature.

**affect the final contribution of the semantic features and does not harm the visual explanations.**

## 5.4.5 | Does DFF adapt to other visual backbones?

DFF is designed to perform concept discovery in CNN-based visual backbones. However, current pre-trained VL models' vision encoders often rely on different architectures, such as ViT (Dosovitskiy et al., 2020b), Faster Recurrent CNN (FRCNN) (Ren et al., 2015b), or their variants. In this section, we show how DFF can be adapted to these architectures. Moreover, we provide an alternative solution to perform model-agnostic semantic feature extraction, which is applicable to any architecture.

**Vision Transformers**   In order to apply DFF to ViT encodings, we need to take into account two substantial differences with respect to CNNs: (1) firstly, ViT splits the image into a grid of patches and generates an embedding vector for each patch. To obtain an activation matrix, each embedding vector is stacked together and a special vector is added in position 0 to indicate the beginning of the sequence. Differently from CNNs, the spatial information related to a patch is lost in the encoding process and added later on, by concatenating a positional embedding to the embedding vectors. (2) Secondly, ViT activations contain both positive and negative values, differently from CNNs which generate only positive activations.

As described in Section 5.3.2.2, DFF requires a non-negative activation matrix as it is based on NMF, therefore in order to address (2) we normalize the ViT features to values between 0 and 1.

As a consequence of (1) above, when we apply DFF to the normalized ViT activations, we obtain binary masks with vertical bands, where each band corresponds to a patch in the image. We used the index of the highlighted vectors in the binary mask to select the patches to be grouped together in the semantic features. In this way, **we ob-**

Figure 5.14: Schematic example of how to generate semantic features with DFF from a FRCNN visual backbone. The index of the highlighted band in the binary mask is used to select the bounding boxes corresponding to objects that compose the input features. However, the bounding boxes highly overlap with each other and cover the majority of the pixels in the image.

**tained feature masks by grouping together semantically related patches**. A schematic example is depicted in Figure 5.13.

**FasterRCNNs**   are often used as feature extractors in VL models (Anderson et al., 2018; Tan and Bansal, 2019; Zhang et al., 2021). They extract feature vectors representing bounding boxes of salient objects identified in the image. Similarly to ViT, the FRCNN's activation matrix is a stack of feature vectors, therefore we can extract semantic features, similarly to the method described in Section 5.4.5. However, FRCNNs tend to extract highly overlapping bounding boxes, which results in massively redundant semantic features. This prevents the features from effectively selecting specific semantic content, as they often result in sharing most of the selected area. A schematic example is shown in Figure 5.14, where although DFF manages to cluster semantically related boxes (like *collar, man, neck, sleeve*), it ends up selecting a large portion of the image in a single input feature.

An excessive amount of overlap among the features affects their capability to identify specific semantic concepts. Thus, we conclude that **DFF can be adapted to FRCNN's features but does not produce the desired results of capturing enough fine-grained semantic concepts to support informative explanations**. In the following subsection, we describe an alternative route towards obtaining semantically meaningful visual regions that can act as features for explaining VL models, in cases where the visual backbone does not permit an application of bottom-up, unsupervised methods such as DFF.

### 5.4.5.1 | Beyond DFF: a model-agnostic semantic feature extraction

As shown in the previous sections:

126

(a) DFF (CNN)          (b) DFF (ViT)          (c) STEGO          (d) Superpixel

Figure 5.15: Direct comparison of explanations generated for the caption "riding a dirt bike" from different visual backbones and methods. The two leftmost explanations (Figures 5.15a and 5.15b) are generated from features extracted using DFF and activations of different visual backbones, namely a CNN (Figure 5.15a) and ViT(Figure 5.15b). Figure 5.15c uses semantic masks extracted by a segmentation model (STEGO) and 5.15d uses superpixel features. All the explanations have comparable compute costs, apart from Figure 5.15c, where only 6 features are used.

- the full potential of DFF is evident when applied to CNN-based models;

- it can be adapted to extract features from ViT models, though these features are less detailed due to the initial discretization of the image into patches operated by the model;

- it does not produce satisfactory results on FRCNN activations, because of the redundancy of the bounding boxes extracted by the model.

In order to address limitations coming from the visual backbone's architecture (e.g. in the case of FRCNNs), we propose to use STEGO (Hamilton et al., 2022)[3] a SOTA segmentation model, to extract semantic feature masks. It is unsupervised, meaning that it does not require ground truth labels. As a consequence, the number of features extracted can not be controlled, though in our experiment we observe that it extracts a small number of semantic masks (usually less than 10). This keeps the Shapley value computation low but could limit the number of semantic concepts captured, differently from DFF where the number of features is a controllable hyperparameter.

The biggest advantage of using an off-the-self segmentation model is that it supports the generation of visual explanations, independently of the visual backbone's architecture. On the other hand, we have the downside of no longer relying on the priors of the visual backbone embedded in the captioning model itself. In other words, by using a

---

[3]At the time of this work, STEGO was a SOTA model for semantic segmentation. However, the approach proposed here is agnostic as to the segmentation model used. For example, Segment Anything (Kirillov et al., 2023), a more recent model proposed after the present experiments were completed could yield better results.

segmentation model we exploit external visual priors which are independent from the VL model we want to explain.

In Figure 5.15 we directly compare the visual explanations generated by all methods, DFF on CNN and ViT (Figures 5.15a and 5.15b), STEGO (Figure 5.15c), and superpixel (Figure 5.15d). All the explanations were generated with similar compute costs, apart from STEGO which uses a smaller amount of features (6). As expected, the explanations generated with STEGO's semantic features are more fine-grained than the others, as the model is trained on the semantic segmentation task. However, they come from an external model and do not necessarily reflect the visual priors of the VL model itself. Nevertheless, this provides a flexible solution to adapt the explanation of VL models with visual priors to any visual backbone. Furthermore, any segmentation model can in principle be used.

## 5.4.6 | Discussion

We have now all the elements needed to answer the question posed at the beginning of this Section, namely:*What are the pros and cons of our method based on visual semantic priors in comparison with standard feature selection methods used in VL, based on superpixels?*.

Exploiting the model's visual priors exposes several significant advantages with respect to standard superpixel features. As shown in Section 5.4.3 input features based on the model's visual priors provide more semantically detailed explanations namely, they succeed in emphasizing salient semantic relevant elements to a higher extent in the image, providing also information regarding the intensity of the area of contribution. The semantic nature of the inputs produces more comprehensive explanations (Section 5.4.3.1) than standard superpixel features at a lower compute cost, thus being also more efficient (Section 5.4.3.2).

However, the introduction of semantic visual features introduces several potential issues that we have thoroughly analysed in this section. From the theoretical point of view, our method requires a relaxation of the feature independence assumption (Section 5.4.4.2) which however, does not compromise the validity of the underlying core method (i.e. KernelSHAP) as we empirically show that non-disjoint features do not significantly affect the visual explanation. In fact, forcing the disjointness of semantic features leads to misleading visual explanations. Similarly, different sizes in the input feature dimension, do not significantly affect the final contribution, as we show (in Section 5.4.4.3). Our method is flexible enough to be adapted to Vision Transformers other than CNNs; however, it adapts with difficulty to FRCNNs (as discussed in Section 5.4.5). To overcome this issue we propose using an off-the-self semantic segmentation model

to extract semantic visual features. In light of our work, which finds its primary motivation in exploiting the model's internal semantic priors, we argue that this solution is not optimal, as it relies on external semantic priors. However, it is a reasonable trade-off that allows us to deal with architectures that do not accommodate DFF to extract such priors.

# 5.5 | Human Evaluation

The experiments in the previous section made direct comparisons between our method and superpixel-based explanations for VL generative models. In this section, we report on an evaluation of human participants aiming to assess the benefits and potential limits of our method for human users.

Evaluating XAI techniques is a notoriously challenging task (e.g. Adebayo et al., 2022; Nauta et al., 2023). Here, we take inspiration from the work of Hoffman et al. (2018) and compare the judgments of participants on three qualities, namely *detail, satisfaction* and *completeness* of explanations generated using the two methods under consideration.

## 5.5.1 | Participants

For the purposes of this study, it is important to source judgments from participants who are knowledgeable about machine learning and explainable AI. Relying on crowdsourcing is a risky strategy, as there is no guarantee that participants will be in a position to evaluate *explanations* rather than, say, the quality of model outputs. We therefore recruited 14 researchers (9 male, 5 female; 9 aged $18-30$, 4 aged $31-40$, 1 aged $41-50$) from our own network. All were researchers in AI-related fields and were familiar with XAI methods. Two of these were senior researchers who obtained their PhD more than 5 years ago; all the others were doctoral students at the time the experiment was run. Six participants were native speakers of English; the remainder are fluent or near-fluent speakers.

## 5.5.2 | Design and materials

We randomly selected 40 images from the HL dataset, for which we generated the corresponding answers to questions. In order to create a more challenging scenario, we framed it into a visual question-answering task, thus for each image, we selected one of the available questions and generate the corresponding caption. Moreover, for each

Figure 5.16: Distribution of the Likert scores obtained in the human evaluation for *detail, completeness* and *satisfaction* for both DFF in (orange) and superpixel (in blue). The lower the score the higher the rating.

image-caption pair, we generated visual explanations using both DFF and superpixel features.

Each participant was shown the question, the generated answer, the original image, and the visual explanation which can be either generated by DFF or by superpixel. In order to counterbalance the experimental materials, we divided images randomly into two groups, and further assigned participants randomly to two groups. We rotated items through a 2 (participant group) × 2 (image group) Latin square, such that participants in any experimental group evaluated all images, but each image was always seen once and evaluated in only one condition (DFF or superpixel).[4] The participants were asked to judge explanations based on their agreement with each of the following statements:

- *Detail*: the areas highlighted in the explanation are detailed enough to understand how the model generated the caption;

- *Completeness*: the highlighted areas cover all the regions relevant to the caption;

- *Satisfaction*: based on the areas highlighted in the explanation I feel that I understand how the system explained makes its decisions.

Responses to each dimension were given on a Likert scale from 1 to 5, where 1 corresponds to the total agreement and 5 to total disagreement. For the full evaluation form see Appendix C.

---

[4]In the end, the experiment was completed by 8 participants in one group, and 6 in the other.

### 5.5.3 | Results

As shown in Figure 5.16, DFF-based explanations (in orange) are considered on par with superpixel-based explanations (in blue) in terms of completeness, but at the same time, they are considered more detailed and more satisfactory for human judges. Thus, the score distributions for detail and satisfaction are skewed towards lower scores (the lower the score the higher the rating).

Although the superpixel and DFF methods differ in the judged level of detail of the explanations, they yield attributions that are similarly located in the input image. This is in part due to the fact that in both cases, we are using the same feature attribution method, namely Kernel SHAP. However, in some cases, we observe a certain degree of divergence in the visual explanation, meaning that the two methods assign opposite attributions to similar regions. In Figure 5.17 we show an example where we generate explanations for the question "Where is the picture taken?" and the generated caption "on a dirty road".

The DFF-based explanation (on the right) broadly assigns a positive attribution to the background of the picture, depicting the road, and negative attributions to the subjects, namely the person and the animals. However, the superpixel-based explanation (on the left) assigns attributions to patches that are, at least partially, in contrast with the DFF-based explanation.

This is probably due to the particular configuration of features selected by both methods, which in some instances might select insufficiently detailed regions, preventing the method from highlighting the semantically relevant areas of the image.

In order to quantify this phenomenon we manually inspected the 40 samples used in the human evaluation. We found that around 10% of the explanations diverged to some extent between the two feature selection methods. We analyzed separately this subsample of divergent explanations. As reported in Table 5.2, the average scores given by experimental participants for this subset are overall slightly worse (higher) than the full results (see supplementary material for details).

Nevertheless, the trends observed in relation to Figure 5.16 for the three evaluation criteria still hold. This suggests that this phenomenon does not significantly affect the participants' judgments, except for a slight drop in the perceived quality of the explanations.

**Impact of caption quality**   In qualitative feedback given by participants, some declared that in some instances, their assessment was affected by the correctness of the caption, which in some cases was considered wrong or partially inaccurate. We quan-

131

Figure 5.17: Comparison of divergent explanations for the question: "Where is the picture taken?" and generated caption: "on a dirty road", obtained from superpixel features (on the left) and DFF features (on the right).

| Type | Metric | Mean | Std | Median |
|------|--------|------|-----|--------|
| SP | completeness | 2.48 | 1.38 | 2.0 |
| | detail | 2.46 | 1.42 | 2.0 |
| | satisfaction | 2.51 | 1.51 | 2.0 |
| DFF | completeness | 2.50 | 1.45 | 2.0 |
| | detail | 2.18 | 1.41 | 2.0 |
| | satisfaction | 2.32 | 1.48 | 2.0 |

Table 5.2: Results of the human evaluation, for superpixel-based (SP) and DFF-based (DFF) visual explanation. We report the mean, the standard deviation (std), and the median of the Likert scores. The lower the score the more positive the rating.

tified the inaccuracy of the caption by computing their lexical and semantic similarity with respect to the reference captions, using respectively, BLEU (Papineni et al., 2002b) and Sentence-Bert (Reimers and Gurevych, 2019). We computed the Pearson correlation (Cohen et al., 2009) between the Likert scores and the lexical and semantic similarity previously computed. As expected, given that 1 is maximum agreement and 5 is minimum, the Likert scores slightly but not significantly negatively correlate with both lexical and semantic similarity ($\rho = -0.023$ for lexical similarity and $\rho = -0.004$ for semantic similarity)[5]. This suggests that despite the fact that participants did note the quality of the captions, this did not significantly affect their judgments of the explanations.

In conclusion, we find that assessing visual explanations is a hard task even for specialists in the field. We observe a relatively low inter-annotator agreement for both groups in the Likert judgments (Krippendorff's $\alpha = 0.23$ (Krippendorff, 2004)). How-

---

[5]Note that since in the Likert score, 1 is the maximum agreement and 5 the minimum, a positive correlation corresponds to a negative $\rho$.

ever, besides possible confounding factors, like inaccuracies of the captions and divergent explanations, the DFF-based explanations are generally perceived as higher quality explanations than superpixel-based ones.

## 5.6 | Summary

In this Chapter, we proposed an explainability framework to bridge the gap between multimodality and explainability in image-to-text generative tasks exploiting textual and visual semantics. Our method is developed around SHAP, as it provides a model-agnostic solution with solid theory and desirable properties. We design our approach to address certain crucial limitations of current approaches

First, SHAP-based methods are rarely employed to explain large models as they are extremely expensive to compute. Our solution is efficient and allows an accurate approximation of the Shapley values.

Second, we overcome the limitations of current token-by-token explanations by proposing sentence-based explanations exploiting semantic textual variations which are also more efficient to compute.

Finally, based on the rationale that a model's generative outputs should be explained with reference to the knowledge encoded by the visual backbone, we proposed an unsupervised method based on DFF to extract semantically informative visual features. Using these features rather than superpixels means that we obtain explanations that are cheaper (insofar as more can be gleaned from fewer features) but also more intuitive, especially when compared to superpixel-based approaches.

We took a self-critical stand on our approach by further studying potential limitations araising from the modifications made to adapt Kernel SHAP to a multimodal setting.

We observed that semantic features extracted with DFF may overlap with each other and therefore may generate non-disjoint features. This directly affects the assumption of feature independence, which is a theoretical requirement for SHAP. However, we study (in Section 5.4.4.2) the extent to which this phenomenon is present and how it affects the outcome of our method finding that it does not significantly affect our visual explanations. Thus, our method features a relaxation of this assumption.

Furthermore, the semantic feature extraction, namely DFF, is designed to extract visual priors from CNN-based models. In Section 5.4.5.1, we showed that this method can successfully adapt to Vision Transformers, but not to FasterRCNNs. To overcome this limitation we proposed to use an off-the-shelf segmentation method to extract semantic

features. This solution supports visual explanation, independently of the visual backbone's architecture. However, in view of the motivation of our work, whose main goal is to exploit the model's visual priors to explain its own predictions, we argue that this solution is not optimal, as it relies on external visual priors (i.e. a third-party semantic segmentation model), though, on the other end offers great flexibility to our framework.

Through a human evaluation (Section 5.5), we showed that using semantic priors improves the perceived quality of the explanation, resulting in more detailed and satisfactory explanations than superpixels though matching the same level of completeness.

We leveraged experts in AI as annotators for our evaluations. However, we are aware that evaluating visual explanations for humans can be a hard task. In particular, the task of evaluating XAI is ambiguous, since evaluators are asked to judge the quality of explanations, which is in principle distinct from the quality of model outputs (that is, one can have a satisfactory explanation of an incorrect or infelicitous output). As the qualitative feedback from our evaluation suggests, keeping output quality and explanation quality separate is not always an easy task and this may influence the evaluation outcomes.

Ultimately, our framework is totally modular and it can co-exist with a wide range of possible configurations for all of its components. For example, it is possible to produce token-by-token explanations and still rely on DFF to extract visual features. The core method, Kernel SHAP, can be replaced with another SHAP-based method, and the visual features can be extracted with one of the proposed methods or with any other method of choice.

In this Chapter, we presented an explainability framework for VL generative models with the main goal to expand and popularize XAI methods for generative multimodal setups. Our framework is model-agnostic and allows to compute efficiently sentence-based visual explanations exploiting the visual priors learned by the model. We discussed benefits and potential limitations providing empirical evidences of the correctness of our method, hoping to provide with a useful contribution for the community and foster further research in this direction.

# Conclusions

## 6.1 | What have we learned?

We summarise what we have been able to learn from this thesis, focusing on the main findings addressing our research questions.

**To what extent do current VL generative models ground high-level information?** To address this question, as anticipated in Section 1.3, we developed a dataset which permits the evaluation of such models in a controlled set of conditions.

In Chapter 1, we argued that VL research has to date been predominantly interested in grounding low-level information, namely objects and entities depicted in the visual content (Hodosh et al., 2013b). Although this is an essential capability to achieve good performance on downstream tasks, it does not tell us much about the capability of VL models to handle high-level descriptions. Such descriptions express high-level information that is often used in human communication (Schank and Abelson, 1975) and form a large part of the language, which is grounded in the perceptual world (Bisk et al., 2020).

To fill this gap we collected a new dataset, the HL Dataset (Cafagna et al., 2023b), aligning existing images and low-level captions with crowd-sourced high-level descriptions. Our dataset leverages two levels of abstraction, namely low- and high-level. We characterise as "low-level" those captions that capture all the information described at the level of objects and entities visible in the image. As discussed in Section 1.2.2 and more in detail in Section 2.3.2, low-level information is already present in the current curated and web-crawled VL datasets. However, in the latter, objectivity is enforced through automatic filtering pipelines which do not guarantee the absolute absence of high-level information in the captions (Van Miltenburg, 2016). By "high-level informa-

tion", on the other hand, we intend all that knowledge not directly present in the visual content, but inferable through interpretation based on common assumptions.

With this conceptual definition in mind, we extended an existing VL dataset, namely COCO (Chen et al., 2015b) with high-level descriptions of *scenes, actions* and *rationales*, by asking annotators to interpret the image based on their assumptions and knowledge of the world. Given the intrinsic subjectivity of high-level captions, we took a step further and collected confidence scores in order to measure the plausibility of these captions with respect to the image. These elements altogether provided a unique combination that can be used to analyze linguistic differences between high- and low-level information as well as find multimodal connections between concepts related to them. This is directly implemented in the analyses presented in Chapter 4, where we focused on the *scene* axis, motivated by the particular interest shared by cognitive sciences and multimodal AI research in studying scene perception in humans and how this could help improving multimodal models.

Moreover, the confidence scores provide an additional dimension to reason upon in terms of personal interpretation and general common sense. These aspects are analysed in Section 3.4.

The HL dataset finds interesting application also in generative tasks. It can be used to generate high-level captions, as shown in Section 3.5 as well as narrative-like captions. In Section 3.6 we presented a case study where we used a combination of high-level captions to generate short "narrative", which describes the scene, action and rationale in tandem. Finally, in Section 3.7 we discussed further potential uses not covered in the scope of this thesis.

Our dataset provides a controlled environment where multiple levels of linguistic abstraction are systematically aligned with images. This alignment enabled us to delve deeper into our investigation, specifically addressing RQ1. Focusing on scene descriptions in Chapter 4, we first analysed the capability of current VL models to handle scene descriptions (Section 4.3) and then we studied the impact of the exposure VL models' representations to scene descriptions in generative settings.

In our zero-shot analysis in Section 4.3, we observed that some, but not all, VL models can handle scene descriptions despite being trained on object-centric textual data. However, this was observed under two conditions: large-scale data pre-training and (ii) an image-sentence alignment modeled through a contrastive pre-training loss. Moreover, when scene-description grounding was successful, the model leveraged object-level visual information, showing the ability to relate typical objects to the scene, similarly to human scene understanding (Vo, 2021; Võ and Wolfe, 2013).

After establishing that the pre-trained VL model can to some extent ground high-

level expression, we followed up our investigation by asking: (i) *Can we adapt pre-trained VL models to generate high-level expressions?* and (ii) *How does direct exposure to such data impact the model?*

To answer these questions, we focused on VL generative settings. We directly addressed (i) in Section 3.5, where we showed that VL generative models can easily be adapted to generate high-level descriptions (i.e. scenes, action, and rationales) through fine-tuning. In Section 4.4, we tackled (ii) by performing an in-depth analysis of the impact of the fine-tuning. We found that the models' representations are rich enough to support scene descriptions; however, fine-tuning helps to bridge the gap with infrequent scene types which might have a less robust representation. The largest impact was observed in the attention mechanism which features a different distribution of the attentional resources. This was confirmed by both qualitative and quantitative experiments. In other words, in order to handle scene descriptions, the model re-distributed the attentional resource more evenly over the visual tokens, suggesting a more "holistic" view of the scene, still retraining the capability to rely on diagnostic object information. We concluded that the model's ability to generate scene-level captions are primarily acquired through a change in the self-attention.

**How can XAI methods be extended to provide a reliable window on model performance with linguistic expressions at different levels?** In Chapter 5, we proposed an explainability framework to bridge the gap between multimodality and explainability in image-to-text generative tasks exploiting visual semantics. Our method is developed around KernelSHAP (Lundberg and Lee, 2017), hence, it is model-agnostic and benefits from all the properties defined in the SHAP framework as well as a solid mathematical definition. We exploited textual semantic representations to allow sentence-based explanations, as opposed to standard token-based explanations. This solution is more suitable for explaining high-level expressions whose meaning often relies on the whole sentence rather than single words. Moreover, as an additional benefit, sentence-based explanations are more efficient than token-based. The semantic representations of the visual modality is central in our framework, as it is a key element to producing efficient as well as meaningful explanations. By exploiting visual priors in the model's visual backbone we reduced the KernelSHAP overall compute cost if compared to traditional superpixel-based explanations.

Efficiency is a critical element to be taken into account to make the SHAP methods realistically applicable in deep learning settings such as VL. On this front, we developed a deterministic approximation version of KernelSHAP which allows a higher degree of efficiency. Our experiments showed that our method provides a good approximation of

the Shapley value even with low-computing capability hardware.

Finally, we ran a human validation to assess the quality of our explanations in comparison with traditional methods. Results showed that exploiting the visual semantics produced more detailed explanations as well as similar in terms of consistency and satisfaction with superpixel-based ones.

## 6.2 | Future Work and Open Questions

Given the results summarised above, we identify some possible directions for future work. Although we bound the notion of high-level descriptions to three aspect of the image, namely *scene, action* and *rationale*, we are aware that an image could be appraised on many high-level aspects such as the temporal (*When was the picture taken? What might happen next/before?*) or the relevance (*What strikes you as interesting in this picture?*). Our notion of "high-level description", could be expanded and refined to accommodate more then two level of abstraction (high- and low-level) opening to further research in this direction.

In Section 4.4 we observe that the capability to generalise to high-level scene descriptions with an object-centric pre-training is linked to two main factors: large-scale pretraining dataset and contrastive loss. Starting from these observations, it would be interesting to investigate further this connection: (i) *are these two factors equally important?*; (ii) *could this be due to some leak of high-level information in large web-crawled datasets?* or *(iii) does the large exposure to object-centric information trigger high-level information generalisation?*. A positive answer to the latter question would be of significant interest, as it would imply a scalable emerging capability to infer high-level information starting from object-level information. Moreover, the analysis performed in Chapter 4 could be expanded to *actions and rationales*.

As discussed in Section 3.7, the HL Dataset could enable potential tasks and use cases that have not been explored in this thesis. One of the most interesting unexplored features of this dataset is the confidence score. As already pointed out, an interesting direction for future work is to use them as a training signal to generate captions with different confidence levels.

In Chapter 5 we present a framework to generate sentence-based explanations for VL generative models exploiting semantic priors. We identify some technical limitations that are addressed empirically. On the evaluation front, the human evaluation could be improved to accommodate a more fine-grained assessment of the explanations and understand in which context these kinds of explanations benefit the user the most.

# 6.3 | Final remarks

Starting from the simple idea that an image can be appraised by at least two levels of abstraction, in this thesis, we argue that previous VL research has been focusing on the linguistic grounding of only one of them, namely low-level. In Chapter 2 we emphasized the importance of high-level information in human communication and the potential benefits from the understanding of how this information can be grounded in VL models. Hence, to fill this gap, in Chapter 3 we introduced the HL dataset, a new crowd-sourced dataset that aligns images with low- and high-level descriptions. This data improves our understanding of the high- and low-level information interplay and how we can combine them to enable new VL tasks and benchmarks. In Chapter 4, we used part of this data to study how VL models handle high-level information and how it impacts their inner workings. Leveraging high-level information poses new challenges to explainability techniques in VL research, where XAI is still lagging behind. In Chapter 5 we introduced a new explainability framework suitable to explain VL models handling high-level as well as low-level information by exploiting visual semantics. Our framework can be used in combination with previous techniques, and bring significant improvement in terms of efficiency, constituting a relevant advancement to XAI applied to VL in generative settings.

Our findings reveal several mechanisms involved in this process and enable us to draw parallels with human cognition. However, we do not claim that VL models are models of the mind. Rather, we demonstrate that VL models establish multimodal connections based on relatively simple mechanisms that rely mainly on correlations and statistics learned from the data. This conclusion raises the question of whether we can use these mechanisms to our advantage to develop more transparent and less data-intensive VL models. In other words, do we really need so much data and so many parameters to learn these multimodal connection? After all, we found that they rely on a few simple mechanisms to model such connections, such as the attention and particular objective functions. The quality of the data seems to be the key to enable these models the capability to ground more complex relationships.

We hope that our contributions will inspire further research on high-level grounding in VL and foster a more general approach to multimodal grounding in the field.

## A.1 | Annotation Details

It is important to note that the instructions, shown in Figure A.1 were always visible to the workers during the data collection.

## A.2 | Annotation Costs

In this section, we report the costs related to the data collection.

**High-level caption collection**  Overall, 1033 participants took part in the caption data collection, they were paid $ 0.04 per item corresponding to the hourly minimum rate in the United Kingdom. In total, the data collection cost $ 1938.

**Confidence Scores collection**  The qualification task for confidence scores led to the recruitment of 53 annotators. We found that this task was harder than the high-level caption annotation in terms of complexity but not in terms of execution time which was indeed shorter. Therefore, in order to encourage good quality annotations, we pay $ 0.04 per item. Considering the time needed to perform the task, this corresponds to 4 times the hourly rate of the minimum wage in the United Kingdom. The qualification task and the data collection cost respectively $ 93 and $ 1938.

## A.3 | Item-based analysis

An item in the HL dataset is an image along with all the high-level captions of all the axes. For instance, Figures A.2 and A.3 show the item-wise *diversity score* and *purity score* distribution respectively, along with their average value across the whole dataset. An

**Instructions**:

You are going to see some pictures. Each picture involves one or more people ('the subject'). You will be asked some questions about the picture

Don't think too much, feel free to give your personal interpretation using your knowledge or common sense.

Try to answer using full English sentences. **If you're not sure what the answer could be, give your best guess.**

Avoid using expressions like "I think" or "I suppose" or "Maybe.

**Do not propose options or possibilities** saying for instance: something "or" something else. **Make your best guess** and state the one you choose.

Write a statement, **don't write a one-word answer**, avoid acronyms or slangs and write a **full sentence**.

1. **Where is the picture taken**: give your best guess about the type of place where the action is happening (for example, "in a ski resort");

2. **What is the subject doing**: Try to describe what the people are doing as concisely as possible. If there is more than one person, try to choose a description that captures what all of them are doing (for example, "They are skiing")

3. **Why is the subject doing it**: here, write your best guess about why the person or persons are doing the action (for example, "They are on a family holiday")

The **What** question and the **Why** question **cannot have the same** answer.

The answers must be **written correctly in English**, check the spell and most importantly **don't forget the subject of the sentence in your answer** (he, she, it, they)

Figure A.1: Final version of the instructions presented to the workers during the collection of the high-level captions. These instructions are always visible to the annotators.

item on the right-hand side of the distribution is systematically more consistent across its axes with respect to the measure considered (*purity* or *diversity*). This information can be combined with confidence scores to perform a more fine-rained sample selection. For example in zero-shot testing, we might want to use a hard sample to test our model with, we can select items with similar lexicons, low-semantic purity, and low confidence scores.

# A.4 | Examples of narrative captions

In Figure A.4 we show examples of narrative caption generations from our fine-tuned baselines.

142

Figure A.2: Item-wise diversity score distribution.



Figure A.3: Item-wise purity score distribution.

BLIP: they are playing football in a soccer field and are spending time together
Gold (T5): They are playing in a stadium they are in a game.



BLIP (FT): he skates in a snowy field and wants to enjoy the ride.
Gold (T5): He is snowboarding in a ski resort and he is on vacation.



ClipClap (FT): They are waiting for a bus to take them to the bus station
Gold (T5): at the bus stops he needs to be taken to his destination.



ClipClap (FT): He is skating on a skateboard in a skate park.
Gold (T5): He is skateboarding at a skatepark for fun.



GIT (FT): they are riding horses in the beach, they want to go on vacation.
Gold (T5): They are riding in a beach, they are in a trip..



GIT (FT): the cat is watching the dog in the kitchen, it is watching television.
Gold (T5): Two cats are watching tv in a living room and wait to be served food.

Figure A.4: Examples of captions generated by the fine-tuned (FT) models and corresponding T5-generated (T5) data on the narrative caption generation task.

## B.1 | Image-sentence alignment details

Results for image-sentence alignment experiments are averages over three separate runs for each model.

We leverage the `transformers`[1] library implementation of LXMERT and standard image-sentence alignment head and the original implementation of CLIP[2].

CLIP and LXMERT are tested on the standard alignment task: given an image and either the correct caption or a random caption, the model needs to determine whether the caption correctly aligns with the image.

We use the publicly available implementation of VisualBERT.[3] The image-sentence alignment setting for this model is somewhat different since alignment is modeled as an extension of the next-sentence prediction task in unimodal BERT. VisualBERT takes an image and a correct caption, together with a second caption, which may be correct or randomly selected. The task is to predict whether the second caption correctly aligns with the image+caption pair.

For all experiments, we truncate textual captions to a maximum length of 50 tokens, following standard practice for such models, including CLIP.

## B.2 | Scene vs entities Examples

---

[1]`github.com/huggingface/transformers`
[2]`https://github.com/openai/CLIP`
[3]`https://github.com/uclanlp/visualbert`

*motorcycle*: 1% *desert*: 99%



*motorcycle*: 96% *forest*: 4%



*couch*: 99% *home*: 1%



*couch*: 22% *home*: 78%



*street*: 88% *bus*: 12%



*street*: 1% *bus*: 99%

Figure B.1: Examples of zero-shot scene vs entity one-to-one comparison performed with CLIP.

# B.3 | Self-attention Details

**Attention beyond Layer 1**   At higher layers the attention converges on the special token [SEP], used to separate the *text + object tags* from the *visual* input, as shown in Fig-

(a) Layer 1                                              (b) Layer 6



(c) Layer 12

Figure B.2: Attention matrices for layers 1, 6, and 12. The attention weights progressively gather on the [SEP] token.

ure 4.8. A similar behaviour has been observed analysing BERT's attention Clark et al. (2019b).

Figure B.2 shows how this pattern becomes more pronounced as we move further across the layers, preventing from observing any kind of input interplay. Although the *text*, *object tags* and *visual* sequences can be of different lengths, the [SEP] token sits always in the same position among the inputs, as the padding is always applied to keep the *text + object tags* sequence of the same length. We believe that this regularity is used by the model as a sort of pivot among the inputs. This can cause a high accumulation of attentional resources by the model.

147

## C.1 | Human Evaluation

The instructions given to participants are shown in Figure C.1. The participant is asked to measure the agreement with three statements related to *detail, completeness* and *satisfaction*. Figure C.2 is an example of the form used by participants to evaluate each item.

Below we will introduce you to some concepts to help you in the evaluation process

1) Our **image captioning** system works by generating **answers to questions**. These answers can be full sentences or even simple words or phrases.

Item
1

You can see an example below:



2) A **visual explanation highlights** the **areas** of the image which **positively or negatively contribute** to generate the **caption**. This is helpful to understand how the system uses the information depicted in the image.

Below you can see the **visual explanation** of the caption **"a man drinking a glass of wine"** , *generated by a* captioning system**.**

Overall,
the **blue** areas **positively contribute** to produce the caption;
the **red** areas **negatively contribute** to produce the caption.

Always refer to the color bar on the right-hand side of the explanation as it gives you a numerical reference to understand the scale of intensity of the color, which might change from an explanation to another,



**3) What are we evaluating?**
We are going to evaluate the following properties by asking to what extent you agree with these statements
- **Detail**: the areas highlighted in the explanation are detailed enough to understand how the model generated the caption
- **Completeness**: the highlighted areas cover all the regions relevant for the caption
- **Satisfaction**: based on the areas highlighted in the explanation I feel that I understand how the system explained makes its decisions

Figure C.1: Instruction presented to the participants of the human evaluation.

Item 1

Question: "where is the picture taken?".
Explaining the answer: " at a skatepark".

REMEMBER:

the blue areas positively contribute to produce the caption;
the red areas negatively contribute to produce the caption.

Always refer to the color bar on the right-hand side of the explanation as it gives
you a numerical reference to understand the scale of intensity of the color, which
might change from an explanation to another.



The areas highlighted in the explanation are detailed enough to undertand how     *
the model generated the caption

I totally agree

1  ◯

2  ◯

3  ◯

4  ◯

5  ◯

I totally disagree

Figure C.2: Example of an item presented to the participants of the human evaluation.
It shows the question, the generated caption, the original image, and the visual expla-
nation. In this Figure, we show the statement related to *detail*.

# References

Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4190–4197, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.385. URL `https://aclanthology.org/2020.acl-main.385`.

Julius Adebayo, Michael Muelly, Hal Abelson, and Been Kim. Post-hoc explanations may be ineffective for detecting unknown spurious correlation. In *Proceedings of the 10th International Conference on Learning Representations (ICLR'22)*, page 13, 2022.

Estelle Aflalo, Meng Du, Shao-Yen Tseng, Yongfei Liu, Chenfei Wu, Nan Duan, and Vasudev Lal. Vl-interpret: An interactive visualization tool for interpreting vision-language transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21406–21415, 2022.

Abien Fred Agarap. Deep learning using rectified linear units (relu). *arXiv preprint arXiv:1803.08375*, 2018.

Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, and Peter Anderson. Nocaps: Novel object captioning at scale. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8948–8957, 2019.

Usman Ahmed, Jerry Chun-Wei Lin, and Gautam Srivastava. Fuzzy explainable attention-based deep active learning on mental-health data. In *2021 IEEE International Conference on Fuzzy Systems (FUZZ-IEEE)*, pages 1–6. IEEE, 2021.

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.

Matt D Anderson, Erich W Graf, James H Elder, Krista A Ehinger, and Wendy J Adams. Category systems for real-world scenes. *Journal of Vision*, 21(2)(8):1–31, 2021.

Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Spice: Semantic propositional image caption evaluation. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part V 14*, pages 382–398. Springer, 2016.

Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018.

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015.

Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.

Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005.

Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021a.

Hangbo Bao, Li Dong, and Furu Wei. BEiT: BERT Pre-Training of Image Transformers. *arXiv*, 2106.08254: 1–16, 2021b.

Hangbo Bao, Wenhui Wang, Li Dong, Qiang Liu, Owais Khan Mohammed, Kriti Aggarwal, Subhojit Som, Songhao Piao, and Furu Wei. Vlmo: Unified vision-language pre-training with mixture-of-modality-experts. *Advances in Neural Information Processing Systems*, 35:32897–32912, 2022.

Lawrence W Barsalou et al. Grounded cognition. *Annual review of psychology*, 59(1):617–645, 2008.

Lisa Beinborn, Teresa Botschen, and Iryna Gurevych. Multimodal grounding for language processing. *arXiv preprint arXiv:1806.06371*, 2018.

Yonatan Belinkov and James Glass. Analysis Methods in Neural Language Processing: A Survey. *Transactions of the Association for Computational Linguistics*, 7:49–72, 2019. doi: 10.1162/tacl\\_a\\_00254.

Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? In *Proceedings of the fourth ACM Conference on Fairness, Accountability, and Transparency (FAccT'21)*, Online, 2021. Association for Computing Machinery.

Irving Biederman, Robert J. Mezzanotte, and Jan C. Rabinowitz. Scene perception: Detecting and judging objects undergoing relational violations. *Cognitive Psychology*, 14(2):143–177, 1982. doi: 10.1016/0010-0285(82)90007-X.

Alexander Binder, Sebastian Bach, Gregoire Montavon, Klaus-Robert Müller, and Wojciech Samek. Layer-wise relevance propagation for deep neural network architectures. In *Information science and applications (ICISA) 2016*, pages 913–922. Springer, 2016.

Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. Multimodal datasets: misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*, 2021.

Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, Nicolas Pinto, and Joseph Turian. Experience grounds language. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8718–8735, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.703. URL `https://aclanthology.org/2020.emnlp-main.703`.

Benedikt Boenninghoff, Steffen Hessler, Dorothea Kolossa, and Robert M Nickel. Explainable authorship verification in social media via attention-based similarity learning. In *2019 IEEE International Conference on Big Data (Big Data)*, pages 36–45. IEEE, 2019.

Leo Breiman. Random forests. *Machine learning*, 45:5–32, 2001.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020a.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. *ArXiv*, 2005.14165, 2020b.

Emanuele Bugliarello, Ryan Cotterell, Naoaki Okazaki, and Desmond Elliott. Multimodal pretraining unmasked: Unifying the vision and language berts. *arXiv preprint arXiv:2011.15124*, 2020.

Emanuele Bugliarello, Ryan Cotterell, Naoaki Okazaki, and Desmond Elliott. Multimodal pretraining unmasked: A meta-analysis and a unified framework of vision-and-language BERTs. *Transactions of the Association for Computational Linguistics*, 9:978–994, 2021a. doi: 10.1162/tacl\\_a\\_00408.

Emanuele Bugliarello, Ryan Cotterell, Naoaki Okazaki, and Desmond Elliott. Multimodal pretraining unmasked: A meta-analysis and a unified framework of vision-and-language berts. *Transactions of the Association for Computational Linguistics*, 9:978–994, 2021b.

Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. Coyo-700m: Image-text pair dataset. `https://github.com/kakaobrain/coyo-dataset`, 2022.

Michele Cafagna, Kees van Deemter, Albert Gatt, et al. What vision-language models 'see' when they see scenes. *ArXiv preprint 2109.07301*, 2021.

Michele Cafagna, Kees van Deemter, and Albert Gatt. Understanding cross-modal interactions in V&L models that generate scene descriptions. In *Proceedings of the Workshop on Unimodal and Multimodal Induction of Linguistic Structures (UM-IoS)*, pages 56–72, Abu Dhabi, United Arab Emirates (Hybrid), December 2022. Association for Computational Linguistics. URL `https://aclanthology.org/2022.umios-1.6`.

Michele Cafagna, Lina M. Rojas-Barahona, Kees van Deemter, and Albert Gatt. Interpreting vision and language generative models with semantic visual priors. *Frontiers in Artificial Intelligence*, 6, 2023a. ISSN 2624-8212. doi: 10.3389/frai.2023.1220476.

155

Michele Cafagna, Kees van Deemter, and Albert Gatt. HL Dataset: Visually-grounded Description of Scenes, Actions and Rationales. In *Proceedings of the 16th International Natural Language Generation Conference (INLG'23)*, Prague, Czech Republic, 2023b.

Jize Cao, Zhe Gan, Yu Cheng, Licheng Yu, Yen-Chun Chen, and Jingjing Liu. Behind the scene: Revealing the secrets of pre-trained vision-and-language models. In *European Conference on Computer Vision*, pages 565–580. Springer, 2020.

Min Cao, Shiping Li, Juntao Li, Liqiang Nie, and Min Zhang. Image-text retrieval: A survey on recent research and development. In Lud De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 5410–5417. International Joint Conferences on Artificial Intelligence Organization, 7 2022. doi: 10.24963/ijcai.2022/759. URL `https://doi.org/10.24963/ijcai.2022/759`. Survey Track.

Yue Cao, Mingsheng Long, Jianmin Wang, and Shichen Liu. Deep visual-semantic quantization for efficient image retrieval. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1328–1337, 2017.

Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3558–3568, 2021.

Hanjie Chen, Guangtao Zheng, and Yangfeng Ji. Generating Hierarchical Explanations on Text Classification via Feature Interaction Detection. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5578–5593, Online, July 2020a. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.494. URL `https://aclanthology.org/2020.acl-main.494`.

Jianfu Chen, Polina Kuznetsova, David Warren, and Choi Yejin. Déjà image-captions: A corpus of expressive descriptions in repetition. pages 504–514, 01 2015a. doi: 10.3115/v1/N15-1053.

Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015b.

Xinlei Chen, Hao Fang, Tsung-yi Lin, Ramakrishna Vedantam, C Lawrence Zitnick, Saurabh Gupta, and Piotr Doll. Microsoft COCO Captions : Data Collection and Evaluation Server. *arXiv preprint 1504.00325*, pages 1–7, 2015c.

Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *ECCV*, 2020b.

Myung Jin Choi, Antonio Torralba, and Alan S Willsky. Context models and out-of-context objects. *Pattern Recognition Letters*, 33(7):853–862, 2012.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.

Kenneth Ward Church and Patrick Hanks. Word association norms, mutual information, and lexicography. *Computational Linguistics*, 16(1):22–29, 1990.

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What does BERT look at? An analysis of BERT's attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy, 2019a. Association for Computational Linguistics. doi: 10.18653/v1/w19-4828. URL `https://www.aclweb.org/anthology/W19-4828`.

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. What does BERT look at? an analysis of BERT's attention. In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy, August 2019b. Association for Computational Linguistics. doi: 10.18653/v1/W19-4828. URL `https://aclanthology.org/W19-4828`.

Israel Cohen, Yiteng Huang, Jingdong Chen, Jacob Benesty, Jacob Benesty, Jingdong Chen, Yiteng Huang, and Israel Cohen. Pearson correlation coefficient. *Noise reduction in speech processing*, pages 1–4, 2009.

Edo Collins, Radhakrishna Achanta, and Sabine Susstrunk. Deep feature factorization for concept discovery. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 336–352, 2018.

Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. Explaining transformer-based image captioning models: An empirical analysis. *AI Communications*, 35(2):111–129, 2022.

Adam Dahlgren Lindström, Johanna Björklund, Suna Bensch, and Frank Drewes. Probing multimodal embeddings for linguistic properties: the visual-semantic case. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 730–744, Barcelona, Spain (Online), December 2020. International Committee on Computational Linguistics. doi: 10.18653/v1/2020.coling-main.64. URL `https://aclanthology.org/2020.coling-main.64`.

Dember, West William N, Jolyon Louis, and William Epstein. "perception". In *Encyclopedia Britannica*. September 2023.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

Thierry Deruyttere, Simon Vandenhende, Dusan Grujicic, Luc Van Gool, and Marie Francine Moens. Talk2car: Taking control of your self-driving car. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2088–2098, 2019.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019a. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL `https://aclanthology.org/N19-1423`.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019b. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL `https://www.aclweb.org/anthology/N19-1423`.

Alexey Dosovitskiy, Lucas Beye, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Un-
terthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and
Neil Houlsby. An image is worth 16x16 words: Transformers for image recognitoin at scale. *arXiv*,
2010.11929, 2020a.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Un-
terthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth
16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020b.

Zi-Yi Dou, Aishwarya Kamath, Zhe Gan, Pengchuan Zhang, Jianfeng Wang, Linjie Li, Zicheng Liu, Ce Liu,
Yann LeCun, Nanyun Peng, et al. Coarse-to-fine vision-language pre-training with fusion in the back-
bone. *Advances in neural information processing systems*, 35:32942–32956, 2022.

Ali Farhadi, Ian Endres, Derek Hoiem, and David Forsyth. Describing objects by their attributes. In *2009
IEEE conference on computer vision and pattern recognition*, pages 1778–1785. IEEE, 2009.

Li Fei-Fei, Robert Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE transactions on
pattern analysis and machine intelligence*, 28(4):594–611, 2006.

Patrick Fernandes, Kayo Yin, Graham Neubig, and André F. T. Martins. Measuring and increasing context
usage in context-aware machine translation. In *Proceedings of the 59th Annual Meeting of the Association for
Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume
1: Long Papers)*, pages 6467–6478, Online, August 2021. Association for Computational Linguistics. doi:
10.18653/v1/2021.acl-long.505. URL `https://aclanthology.org/2021.acl-long.505`.

Francis Ferraro, Nasrin Mostafazadeh, Lucy Vanderwende, Jacob Devlin, Michel Galley, Margaret Mitchell,
et al. A survey of current datasets for vision and language research. *arXiv preprint arXiv:1506.06833*, 2015.

Adam Fisch, Kenton Lee, Ming Wei Chang, Jonathan H. Clark, and Regina Barzilay. CAPWAP: Captioning
with a Purpose. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing
(EMNLP)*, Online, 2020a. Association for Computational Linguistics.

Adam Fisch, Kenton Lee, Ming-Wei Chang, Jonathan H Clark, and Regina Barzilay. Capwap: Captioning
with a purpose. *arXiv preprint arXiv:2011.04264*, 2020b.

Aaron Fisher, Cynthia Rudin, and Francesca Dominici. All models are wrong, but many are useful: Learn-
ing a variable's importance by studying an entire class of prediction models simultaneously. *J. Mach.
Learn. Res.*, 20(177):1–81, 2019.

Stella Frank, Emanuele Bugliarello, and Desmond Elliott. Vision-and-Language or Vision-for-Language?
On Cross-Modal Influence in Multimodal Transformers. *ArXiv*, 2109.04448, 2021a.

Stella Frank, Emanuele Bugliarello, and Desmond Elliott. Vision-and-language or vision-for-language?
on cross-modal influence in multimodal transformers. In *Proceedings of the 2021 Conference on Empirical
Methods in Natural Language Processing*, pages 9847–9857, Online and Punta Cana, Dominican Republic,
November 2021b. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.775.
URL `https://aclanthology.org/2021.emnlp-main.775`.

Stanislav Frolov, Tobias Hinz, Federico Raue, Jörn Hees, and Andreas Dengel. Adversarial text-to-image
synthesis: A review. *Neural Networks*, 144:187–209, 2021.

Zhe Gan, Linjie Li, Chunyuan Li, Lijuan Wang, Zicheng Liu, Jianfeng Gao, et al. Vision-language pre-training: Basics, recent advances, and future trends. *Foundations and Trends® in Computer Graphics and Vision*, 14(3–4):163–352, 2022.

Albert Gatt and Emiel Krahmer. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. *Journal of Artificial Intelligence Research*, 61:65–170, 2018.

Yash Goyal, Tejas Khot, Aishwarya Agrawal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR'17)*, pages 6904–6913, 2017. doi: 10.1007/s11263-018-1116-0.

Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. Annotation artifacts in natural language inference data. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2017. URL https://aclanthology.org/N18-2017.

Mark Hamilton, Zhoutong Zhang, Bharath Hariharan, Noah Snavely, and William T Freeman. Unsupervised semantic segmentation by distilling feature correspondences. *arXiv preprint arXiv:2203.08414*, 2022.

Stevan Harnad. The symbol grounding problem. *Physica*, D42(1990):335–346, 1990.

Andrew F Hayes and Klaus Krippendorff. Answering the call for a standard reliability measure for coding data. *Communication methods and measures*, 1(1):77–89, 2007.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. corr abs/1512.03385 (2015), 2015a.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1026–1034, 2015b. doi: 10.1109/ICCV.2015.123.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

Lisa Anne Hendricks and Aida Nematzadeh. Probing Image-Language Transformers for Verb Understanding. In *Findings ofthe Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3635–3644, Online, 2021a. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.318.

Lisa Anne Hendricks and Aida Nematzadeh. Probing image-language transformers for verb understanding. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3635–3644, Online, August 2021b. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-acl.318. URL https://aclanthology.org/2021.findings-acl.318.

Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. Generating Visual Explanations. In *Proceedings of the 2016 European Conference on Computer Vision (ECCV'16)*, Amsterdam, 2016. URL http://arxiv.org/abs/1603.08507. arXiv: 1603.08507.

159

Lisa Anne Hendricks, Ronghang Hu, Trevor Darrell, and Zeynep Akata. Generating counterfactual explanations with natural language. *arXiv preprint arXiv:1806.09809*, 2018.

Lisa Anne Hendricks, John Mellor, Rosalia Schneider, Jean-Baptiste Alayrac, and Aida Nematzadeh. Decoupling the role of data, attention, and losses in multimodal transformers. *Transactions of the Association for Computational Linguistics*, 9:570–585, 2021.

Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.

Micah Hodosh, Peter Young, and Julia Hockenmaier. Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, 47:853–899, 2013a. ISSN 10769757. doi: 10.1613/jair.3994.

Micah Hodosh, Peter Young, and Julia Hockenmaier. Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, 47:853–899, 2013b.

Micah Hodosh, Peter Young, and Julia Hockenmaier. Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, 47:853–899, 2013c.

Robert R Hoffman, Shane T Mueller, Gary Klein, and Jordan Litman. Metrics for explainable ai: Challenges and prospects. *arXiv preprint arXiv:1812.04608*, 2018.

Xudong Hong, Asad Sayeed, Khushboo Mehra, Vera Demberg, and Bernt Schiele. Visual writing prompts: Character-grounded story generation with curated image sequences. *Transactions of the Association for Computational Linguistics*, 11:565–581, 2023.

MD Zakir Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. A comprehensive survey of deep learning for image captioning. *ACM Computing Surveys (CsUR)*, 51(6):1–36, 2019.

David M. Howcroft, Anya Belz, Miruna-Adriana Clinciu, Dimitra Gkatzia, Sadid A. Hasan, Saad Mahamood, Simon Mille, Emiel van Miltenburg, Sashank Santhanam, and Verena Rieser. Twenty years of confusion in human evaluation: NLG needs evaluation sheets and standardised definitions. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 169–182, Dublin, Ireland, December 2020. Association for Computational Linguistics. URL `https://aclanthology.org/2020.inlg-1.23`.

Junjie Hu, Yu Cheng, Zhe Gan, Jingjing Liu, Jianfeng Gao, and Graham Neubig. What makes a good story? designing composite rewards for visual storytelling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7969–7976, 2020.

Xiaowei Hu, Xi Yin, Kevin Lin, Lei Zhang, Jianfeng Gao, Lijuan Wang, and Zicheng Liu. Vivo: Visual vocabulary pre-training for novel object captioning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(2):1575–1583, May 2021a.

Xiaowei Hu, Xi Yin, Kevin Lin, Lei Zhang, Jianfeng Gao, Lijuan Wang, and Zicheng Liu. Vivo: Visual vocabulary pre-training for novel object captioning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1575–1583, 2021b.

160

Xiaowei Hu, Zhe Gan, Jianfeng Wang, Zhengyuan Yang, Zicheng Liu, Yumao Lu, and Lijuan Wang. Scaling up vision-language pre-training for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 17980–17989, June 2022a.

Xiaowei Hu, Zhe Gan, Jianfeng Wang, Zhengyuan Yang, Zicheng Liu, Yumao Lu, and Lijuan Wang. Scaling up vision-language pre-training for image captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 17980–17989, 2022b.

Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. Visual storytelling. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 1233–1239, 2016.

Zhicheng Huang, Zhaoyang Zeng, Bei Liu, Dongmei Fu, and Jianlong Fu. Pixel-bert: Aligning image pixels with text by deep multi-modal transformers. *arXiv preprint arXiv:2004.00849*, 2020.

Zhicheng Huang, Zhaoyang Zeng, Yupan Huang, Bei Liu, Dongmei Fu, and Jianlong Fu. Seeing out of the box: End-to-end pre-training for vision-language representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12976–12985, 2021.

Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019.

Yuqi Huo, Manli Zhang, Guangzhen Liu, Haoyu Lu, Yizhao Gao, Guoxing Yang, Jingyuan Wen, Heng Zhang, Baogui Xu, Weihao Zheng, et al. Wenlan: Bridging vision and language by large-scale multi-modal pre-training. *arXiv preprint arXiv:2103.06561*, 2021.

Gabriel Ilharco, Rowan Zellers, Ali Farhadi, and Hannaneh Hajishirzi. Probing Contextual Language Models for Common Ground with Visual Representations. *arXiv*, 2005.00619, 2020.

L Itti and C Koch. Computational modelling of visual attention. *Nature reviews neuroscience*, 2(3):194–203, 2001. doi: 10.1038/35058500.

Sai Muralidhar Jayanthi, Danish Pruthi, and Graham Neubig. Neuspell: A neural spelling correction toolkit. *arXiv preprint arXiv:2010.11085*, 2020.

Jeya Vikranth Jeyakumar, Joseph Noor, Yu-Hsi Cheng, Luis Garcia, and Mani Srivastava. How can i explain this to you? an empirical study of deep neural network explanation methods. *Advances in Neural Information Processing Systems*, 33:4211–4222, 2020.

Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pages 4904–4916. PMLR, 2021.

Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2901–2910, 2017.

161

Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1780–1790, 2021.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.

Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014.

Ilker Kesen, Andrea Pedrotti, Mustafa Dogan, Michele Cafagna, Emre Can Acikgoz, Letitia Parcalabescu, Iacer Calixto, Anette Frank, Albert Gatt, Aykut Erdem, et al. Vilma: A zero-shot benchmark for linguistic and temporal grounding in video-language models. *arXiv preprint arXiv:2311.07022*, 2023.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.

Enja Kokalj, Blaž Škrlj, Nada Lavrač, Senja Pollak, and Marko Robnik-Šikonja. Bert meets shapley: Extending shap explanations to transformer-based classifiers. In *Proceedings of the EACL Hackashop on News Media Content Analysis and Automated Report Generation*, pages 16–21, 2021.

Klaus Krippendorff. Reliability in content analysis: Some common misconceptions and recommendations. *Human communication research*, 30(3):411–433, 2004.

Klaus Krippendorff. *Content analysis: An introduction to its methodology*. Sage publications, 2018.

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *Int. J. Comput. Vision*, 123(1):32–73, may 2017a. ISSN 0920-5691. doi: 10.1007/s11263-016-0981-7.

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123(1):32–73, 2017b.

Taku Kudo and John Richardson. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*, 2018.

Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Alexander Kolesnikov, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International Journal of Computer Vision*, 128(7):1956–1981, 2020.

Jing Lei, Max G'Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):1094–1111, 2018.

V. I. Levenshtein. Binary Codes Capable of Correcting Deletions, Insertions and Reversals. *Soviet Physics Doklady*, 10:707, February 1966.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.

Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. SEED-Bench: Benchmarking Multimodal LLMs with Generative Comprehension, July 2023a. URL `http://arxiv.org/abs/2307.16125`. arXiv:2307.16125 [cs].

Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip Torr. Controllable text-to-image generation. *Advances in Neural Information Processing Systems*, 32, 2019.

Chenliang Li, Ming Yan, Haiyang Xu, Fuli Luo, Wei Wang, Bin Bi, and Songfang Huang. SemVLP: Vision-Language Pre-training by Aligning Semantics at Multiple Levels. *arXiv preprint*, 2103.07829, 2021a.

Chenliang Li, Haiyang Xu, Junfeng Tian, Wei Wang, Ming Yan, Bin Bi, Jiabo Ye, Hehong Chen, Guohai Xu, Zheng Cao, et al. mplug: Effective and efficient vision-language learning by cross-modal skip-connections. *arXiv preprint arXiv:2205.12005*, 2022a.

Feng Li, Hao Zhang, Yi-Fan Zhang, Shilong Liu, Jian Guo, Lionel M Ni, PengChuan Zhang, and Lei Zhang. Vision-language intelligence: Tasks, representation learning, and large models. *arXiv preprint arXiv:2203.01922*, 2022b.

Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Daxin Jiang. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 11336–11344, 2020a.

Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021b.

Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022c.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023b.

Liangzhi Li, Bowen Wang, Manisha Verma, Yuta Nakashima, Ryo Kawasaki, and Hajime Nagahara. Scouter: Slot attention-based classifier for explainable image recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1046–1055, 2021c.

Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. What Does BERT with Vision Look At? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL'20)*, pages 5265–5275, Online, 2020b. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.469.

Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. What does BERT with vision look at? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5265–5275, Online, July 2020c. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.469. URL `https://aclanthology.org/2020.acl-main.469`.

Wei Li, Can Gao, Guocheng Niu, Xinyan Xiao, Hao Liu, Jiachen Liu, Hua Wu, and Haifeng Wang. Unimo: Towards unified-modal understanding and generation via cross-modal contrastive learning. *arXiv preprint arXiv:2012.15409*, 2020d.

Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*, pages 121–137. Springer, 2020e.

Yixin Li, Chen Li, Xiaoyan Li, Kai Wang, Md Mamunur Rahaman, Changhao Sun, Hao Chen, Xinran Wu, Hong Zhang, and Qian Wang. A comprehensive review of markov random field and conditional random field approaches in pathology image analysis. *Archives of Computational Methods in Engineering*, 29(1):609–639, 2022d.

Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004a. Association for Computational Linguistics. URL `https://aclanthology.org/W04-1013`.

Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004b.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 740–755, Cham, 2014a. Springer International Publishing. ISBN 978-3-319-10602-1.

Tsung Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. In *Proceedings of the 2014 European Conference on Comupter Vision (ECCV'14)*, volume 8693 LNCS, pages 740–755, Berlin and Heidelberg, 2014b. Springer. ISBN 978-3-319-10601-4. doi: 10.1007/978-3-319-10602-1\\_48.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014c.

Fangyu Liu, Emanuele Bugliarello, Edoardo Maria Ponti, Siva Reddy, Nigel Collier, and Desmond Elliott. Visually grounded reasoning across languages and cultures. *arXiv preprint arXiv:2109.13238*, 2021a.

Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial reasoning. *Transactions of the Association for Computational Linguistics*, 11:635–651, 2023.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021b.

Siqu Long, Feiqi Cao, Soyeon Caren Han, and Haiqin Yang. Vision-and-language pretrained models: A survey. *arXiv preprint arXiv:2204.07356*, 2022.

Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

Octavio Loyola-Gonzalez. Black-box vs. white-box: Understanding their advantages and weaknesses from a practical point of view. *IEEE access*, 7:154096–154113, 2019.

Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32, 2019.

Stephanie Lukin, Reginald Hobbs, and Clare Voss. A pipeline for creative visual storytelling. In *Proceedings of the First Workshop on Storytelling*, pages 20–32, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-1503. URL `https://aclanthology.org/W18-1503`.

Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30, 2017.

Hao Ma, Jianke Zhu, Michael Rung Tsong Lyu, and Irwin King. Bridging the semantic gap between image contents and tags. *IEEE Transactions on Multimedia*, 12(5):462–473, 2010. ISSN 15209210. doi: 10.1109/TMM.2010.2051360.

Kiwan Maeng, Alexei Colin, and Brandon Lucia. Alpaca: Intermittent execution without checkpoints. *Proceedings of the ACM on Programming Languages*, 1(OOPSLA):1–30, 2017.

George L. Malcolm, Iris I.A. Groen, and Chris I. Baker. Making Sense of Real-World Scenes. *Trends in Cognitive Sciences*, 20(11):843–856, 2016. doi: 10.1016/j.tics.2016.09.003.

Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 11–20, 2016.

Tom McCoy, Ellie Pavlick, and Tal Linzen. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3428–3448, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1334. URL `https://aclanthology.org/P19-1334`.

Drew McDermott. Minds, brains, programs, and persons. *Behavioral and Brain Sciences*, 5:339 – 341, 1982.

Tomas Mikolov, G.s Corrado, Kai Chen, and Jeffrey Dean. Efficient estimation of word representations in vector space. pages 1–12, 01 2013a.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013b.

165

Ron Mokady, Amir Hertz, and Amit H Bermano. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021.

Christoph Molnar. *Interpretable machine learning*. Lulu. com, 2020.

Edoardo Mosca, Ferenc Szigeti, Stella Tragianni, Daniel Gallagher, and Georg Groh. SHAP-based explanation methods: A review for NLP interpretability. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 4593–4603, Gyeongju, Republic of Korea, October 2022a. International Committee on Computational Linguistics. URL `https://aclanthology.org/2022.coling-1.406`.

Edoardo Mosca, Ferenc Szigeti, Stella Tragianni, Daniel Gallagher, and Georg Groh. Shap-based explanation methods: A review for nlp interpretability. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 4593–4603, 2022b.

W James Murdoch, Chandan Singh, Karl Kumbier, Reza Abbasi-Asl, and Bin Yu. Interpretable machine learning: definitions, methods, and applications. *arXiv preprint arXiv:1901.04592*, 2019.

Meike Nauta, Jan Trienes, Shreyasi Pathak, Elisa Nguyen, Michelle Peters, Yasmin Schmitt, Jörg Schlötterer, Maurice van Keulen, and Christin Seifert. From Anecdotal Evidence to Quantitative Evaluation Methods: A Systematic Review on Evaluating Explainable AI. *ACM Computing Surveys*, February 2023. ISSN 0360-0300. doi: 10.1145/3583558.

Joe O'Connor and Jacob Andreas. What context features can transformer language models use? In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 851–864, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.70. URL `https://aclanthology.org/2021.acl-long.70`.

Aude Oliva and Antonio Torralba. The role of context in object recognition. *Trends in cognitive sciences*, 11 (12):520–527, 2007. ISSN 1364-6613. doi: 10.1016/j.tics.2007.09.009.

Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. In *Proceedings of the 2011 Conference on Advances in Neural Information Processing Systems (NIPS'11)*, pages 1143–1151, Granada, Spain, 2011a. Curran Associates Ltd. URL `http://machinelearning.wustl.edu/mlpapers/paper\_files/NIPS2011\_0671.pdf`.

Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. *Advances in neural information processing systems*, 24, 2011b.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, July 2002a. Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL `https://aclanthology.org/P02-1040`.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evalua-
tion of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational
Linguistics*, pages 311–318, 2002b.

Letitia Parcabalescu, Albert Gatt, Annette Frank, and Iacer Calixto. Seeing past words: Testing the cross-
modal capabilities of pretrained v&l models on counting tasks. In *Proceedings of the Workshop Beyond
Language: Multimodal Semantic Representations (MMSR'21)*, Groningen, The Netherlands, 2021. URL
`https://arxiv.org/abs/2012.12352`.

Letitia Parcalabescu and Anette Frank. Mm-shap: A performance-agnostic metric for measuring multi-
modal contributions in vision and language models & tasks. *arXiv preprint arXiv:2212.08158*, 2022.

Letitia Parcalabescu and Anette Frank. MM-SHAP: A performance-agnostic metric for measuring mul-
timodal contributions in vision and language models & tasks. In *Proceedings of the 61st Annual Meet-
ing of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4032–4059, Toronto,
Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.223. URL
`https://aclanthology.org/2023.acl-long.223`.

Letitia Parcalabescu, Albert Gatt, Anette Frank, and Iacer Calixto. Seeing past words: Testing the cross-
modal capabilities of pretrained v&l models on counting tasks. *arXiv preprint arXiv:2012.12352*, 2020.

Letitia Parcalabescu, Michele Cafagna, Lilitta Muradjan, Anette Frank, Iacer Calixto, and Albert Gatt.
VALSE: A task-independent benchmark for vision and language models centered on linguistic phenom-
ena. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1:
Long Papers)*, pages 8253–8280, Dublin, Ireland, May 2022a. Association for Computational Linguistics.
doi: 10.18653/v1/2022.acl-long.567. URL `https://aclanthology.org/2022.acl-long.567`.

Letitia Parcalabescu, Michele Cafagna, Lilitta Muradjan, Anette Frank, Iacer Calixto, and Albert Gatt.
VALSE: A task-independent benchmark for vision and language models centered on linguistic phenom-
ena. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1:
Long Papers)*, pages 8253–8280, Dublin, Ireland, May 2022b. Association for Computational Linguistics.
doi: 10.18653/v1/2022.acl-long.567. URL `https://aclanthology.org/2022.acl-long.567`.

Letitia Parcalabescu, Michele Cafagna, Lilitta Muradjan, Anette Frank, Iacer Calixto, and Albert Gatt.
VALSE: A task-independent benchmark for vision and language models centered on linguistic phenom-
ena. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1:
Long Papers)*, pages 8253–8280, Dublin, Ireland, May 2022c. Association for Computational Linguistics.
doi: 10.18653/v1/2022.acl-long.567. URL `https://aclanthology.org/2022.acl-long.567`.

Jae Sung Park, Chandra Bhagavatula, Roozbeh Mottaghi, Ali Farhadi, and Yejin Choi. Visualcomet: Rea-
soning about the dynamic context of a still image. In *European Conference on Computer Vision*, pages
508–524. Springer, 2020.

Gustavo Penha and Claudia Hauff. What does bert know about books, movies and music? probing bert
for conversational recommendation. In *Proceedings of the 14th ACM Conference on Recommender Systems*,
pages 388–397, 2020.

Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1162. URL `https://aclanthology.org/D14-1162`.

Vitali Petsiuk, Abir Das, and Kate Saenko. Rise: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421*, 2018.

Sandro Pezzelle, Claudio Greco, Greta Gandolfi, Eleonora Gualdoni, and Raffaella Bernardi. Be Different to Be Better! A Benchmark to Leverage the Complementarity of Language and Vision. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2751–2767, Online, 2020a. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.248.

Sandro Pezzelle, Claudio Greco, Greta Gandolfi, Eleonora Gualdoni, and Raffaella Bernardi. Be different to be better! a benchmark to leverage the complementarity of language and vision. In *Findings of the association for computational linguistics: EMNLP 2020*, pages 2751–2767, 2020b.

Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pages 2641–2649, 2015.

Jordi Pont-Tuset, Jasper Uijlings, Soravit Changpinyo, Radu Soricut, and Vittorio Ferrari. Connecting Vision and Language with Localized Narratives. *arXiv*, 1912.03098, 2019.

Jordi Pont-Tuset, Jasper Uijlings, Soravit Changpinyo, Radu Soricut, and Vittorio Ferrari. Connecting vision and language with localized narratives. In *European Conference on Computer Vision*, pages 647–664. Springer, 2020.

Matt Post. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels, October 2018. Association for Computational Linguistics. URL `https://www.aclweb.org/anthology/W18-6319`.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning Transferable Visual Models From Natural Language Supervision. *arXiv preprint*, 2103.00020, 2021a.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021b.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.

Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.

Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.

Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In *Advances in Neural Information Processing Systems 28 (NeurIPS 2015)*, Montreal, Canada, 2015a.

Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015b.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Model-agnostic interpretability of machine learning. *arXiv preprint arXiv:1606.05386*, 2016a.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016b.

Robin Rombach, Patrick Esser, and Björn Ommer. Making sense of cnns: Interpreting deep representations and their invariances with inns. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVII 16*, pages 647–664. Springer, 2020.

Cynthia Rudin, Chaofan Chen, Zhi Chen, Haiyang Huang, Lesia Semenova, and Chudi Zhong. Interpretable machine learning: Fundamental principles and 10 grand challenges. *Statistic Surveys*, 16:1–85, 2022.

Ananya B Sai, Akash Kumar Mohankumar, and Mitesh M Khapra. A survey of evaluation metrics used for nlg systems. *ACM Computing Surveys (CSUR)*, 55(2):1–39, 2022.

Emmanuelle Salin, Badreddine Farah, Stéphane Ayache, and Benoit Favre. Are vision-language transformers learning multimodal representations? a probing perspective. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11248–11257, 2022.

Roger C. Schank and Robert P. Abelson. Scripts, plans, and knowledge. In *Proceedings of the 4th International Joint Conference on Artificial Intelligence - Volume 1*, IJCAI'75, page 151–157, San Francisco, CA, USA, 1975. Morgan Kaufmann Publishers Inc.

Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.

Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022.

John R Searle. *Minds, brains and science*. Harvard university press, 1984.

169

Julie S. Self, Jamie Siegart, Munashe Machoko, Enton Lam, and Michelle R Greene. Diagnostic Objects Contribute to Late – But Not Early– Visual Scene Processing. *Journal of Vision*, 19:227, 2019.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.704. URL `https://aclanthology.org/2020.acl-main.704`.

Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.

Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8430–8439, 2019.

Lloyd S Shapley et al. A value for n-person games. 1953.

Himanshu Sharma, Manmohan Agrahari, Sujeet Kumar Singh, Mohd Firoj, and Ravi Kumar Mishra. Image captioning: a comprehensive survey. In *2020 International Conference on Power Electronics & IoT Applications in Renewable Energy and its Control (PARC)*, pages 325–328. IEEE, 2020.

Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL'18)*, pages 2556–2565, Melbourne, Australia, 2018a. Association for Computational Linguistics. doi: 10.18653/v1/p18-1238.

Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, Melbourne, Australia, July 2018b. Association for Computational Linguistics. doi: 10.18653/v1/P18-1238. URL `https://aclanthology.org/P18-1238`.

Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of ACL*, 2018c.

Ravi Shekhar, Sandro Pezzelle, Yauhen Klimovich, Aurélie Herbelot, Moin Nabi, Enver Sangineto, and Raffaella Bernardi. Foil it! find one mismatch between image and language caption. *arXiv preprint arXiv:1705.01359*, 2017.

Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black box: Learning important features through propagating activation differences. *arXiv preprint arXiv:1605.01713*, 2016.

Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *International conference on machine learning*, pages 3145–3153. PMLR, 2017.

Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.

170

Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15638–15650, 2022.

Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.

Krishna Srinivasan, Karthik Raman, Jiecao Chen, Michael Bendersky, and Marc Najork. Wit: Wikipedia-based image text dataset for multimodal multilingual machine learning. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '21, page 2443–2449, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450380379. doi: 10.1145/3404835.3463257. URL https://doi.org/10.1145/3404835.3463257.

Yash Srivastava, Vaishnav Murali, Shiv Ram Dubey, and Snehasis Mukherjee. Visual question answering using deep learning: A survey and performance analysis. In *Computer Vision and Image Processing: 5th International Conference, CVIP 2020, Prayagraj, India, December 4-6, 2020, Revised Selected Papers, Part II 5*, pages 75–86. Springer, 2021.

Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. A Corpus of Natural Language for Visual Reasoning. In *Proceedings ofthe 55th Annual Meeting ofthe Association for Computational Linguistics (ACL'17)*, pages 217–223, Vancouver, BC, 2017a. Association for Computational Linguistics. doi: 10.18653/v1/P17-2034. URL https://doi.org/10.18653/v1/P17-2034.

Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. A corpus of natural language for visual reasoning. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 217–223, Vancouver, Canada, July 2017b. Association for Computational Linguistics. doi: 10.18653/v1/P17-2034. URL https://aclanthology.org/P17-2034.

Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. *arXiv preprint arXiv:1811.00491*, 2018.

Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6418–6428, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1644. URL https://aclanthology.org/P19-1644.

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR, 2017.

Hao Tan and Mohit Bansal. LXMERT: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5100–5111, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1514. URL https://aclanthology.org/D19-1514.

Marc Tanti, Albert Gatt, and Kenneth P Camilleri. Where to put the image in an image caption generator. *Natural Language Engineering*, 24(3):467–489, 2018.

Jacopo Teneggi, Alexandre Luster, and Jeremias Sulam. Fast Hierarchical Games for Image Explanations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–11, 2022. ISSN 0162-8828, 2160-9292, 1939-3539. doi: 10.1109/TPAMI.2022.3189849. arXiv:2104.06164 [cs].

Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5238–5248, 2022.

Antonio Torralba, Aude Oliva, Monica S. Castelhano, and John M. Henderson. Contextual guidance of eye movements and attention in real-world scenes: The role of global features in object search. *Psychological Review*, 113(4):766–786, 2006. doi: 10.1037/0033-295X.113.4.766.

Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii: Revenge of the vit. In *European Conference on Computer Vision*, pages 516–533. Springer, 2022.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

Shikhar Tuli, Ishita Dasgupta, Erin Grant, and Thomas L. Griffiths. Are Convolutional Neural Networks or Transformers more like human vision? *arXiv*, 2105.07197, 2021.

Tomoki Uchiyama, Naoya Sogi, Koichiro Niinuma, and Kazuhiro Fukui. Visually explaining 3d-cnn predictions for video classification with an adaptive occlusion sensitivity analysis. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1513–1522, 2023.

Emiel Van Miltenburg. Stereotyping and bias in the flickr30k dataset. *arXiv preprint arXiv:1605.06083*, 2016.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017a.

Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017b.

Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575, 2015.

Ramakrishna Vedantam, Arthur Szlam, Maximillian Nickel, Ari Morcos, and Brenden M Lake. Curi: A benchmark for productive concept learning under uncertainty. In *International Conference on Machine Learning*, pages 10519–10529. PMLR, 2021.

Melissa Le Hoa Vo. The meaning and structure of scenes. *Vision Research*, 181:10–20, 2021. doi: 10.1016/j.visres.2020.11.003.

Hoa Trong Vu, Claudio Greco, Aliia Erofeeva, Somayeh Jafaritazehjan, Guido Linders, Marc Tanti, Alberto Testoni, Raffaella Bernardi, and Albert Gatt. Grounded textual entailment. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2354–2368, Santa Fe, New Mexico, USA, August 2018. Association for Computational Linguistics. URL `https://aclanthology.org/C18-1199`.

Melissa L-H Võ and Jeremy M Wolfe. Differential electrophysiological signatures of semantic and syntactic scene processing. *Psychological science*, 24(9):1816–1823, 2013.

Melissa Le-Hoa Võ. The meaning and structure of scenes. *Vision Research*, 181:10–20, 2021. ISSN 0042-6989. doi: https://doi.org/10.1016/j.visres.2020.11.003.

Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.

Jonas Wallat, Fabian Beringer, Abhijit Anand, and Avishek Anand. Probing bert for ranking abilities. In *European Conference on Information Retrieval*, pages 255–273. Springer, 2023.

Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. Git: A generative image-to-text transformer for vision and language. *arXiv preprint arXiv:2205.14100*, 2022a.

Jing Wang, Yingwei Pan, Ting Yao, Jinhui Tang, and Tao Mei. Convolutional auto-encoding of sentence topics for image paragraph generation. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, IJCAI'19, page 940–946. AAAI Press, 2019. ISBN 9780999241141.

Josiah Wang, Pranava Swaroop Madhyastha, and Lucia Specia. Object counts! bringing explicit detections back into image captioning. In Marilyn Walker, Heng Ji, and Amanda Stent, editors, *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2180–2193, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1198. URL `https://aclanthology.org/N18-1198`.

Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pages 23318–23340. PMLR, 2022b.

Zhecan Wang, Haoxuan You, Yicheng He, Wenhao Li, Kai-Wei Chang, and Shih-Fu Chang. Understanding ME? multimodal evaluation for fine-grained visual commonsense. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9212–9224, Abu Dhabi, United Arab Emirates, December 2022c. Association for Computational Linguistics. URL `https://aclanthology.org/2022.emnlp-main.626`.

Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. Simvlm: Simple visual language model pretraining with weak supervision. *arXiv preprint arXiv:2108.10904*, 2021.

William Webber, Alistair Moffat, and Justin Zobel. A similarity measure for indefinite rankings. *ACM Transactions on Information Systems (TOIS)*, 28(4):1–38, 2010.

Haiyang Wei, Zhixin Li, Feicheng Huang, Canlong Zhang, Huifang Ma, and Zhongzhi Shi. Integrating scene semantic knowledge into image captioning. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 17(2):1–22, 2021.

Junda Wu, Tong Yu, and Shuai Li. Deconfounded and explainable interactive vision-language retrieval of complex scenes. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 2103–2111, 2021.

Jianxiong Xiao, James Hays, Krista A. Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'10)*, pages 3485–3492, 2010. doi: 10.1109/CVPR.2010.5539970.

Ning Xie, Farley Lai, Derek Doran, and Asim Kadav. Visual entailment: A novel task for fine-grained image understanding. *arXiv preprint arXiv:1901.06706*, 2019.

Xu Yang, Hanwang Zhang, Guojun Qi, and Jianfei Cai. Causal attention for vision-language tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9847–9857, 2021.

Yu Yang, Besmira Nushi, Hamid Palangi, and Baharan Mirzasoleiman. Mitigating spurious correlations in multi-modal models during fine-tuning. *arXiv preprint arXiv:2304.03916*, 2023.

Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014a. doi: 10.1162/tacl\_a\_00166.

Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014b.

Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022.

Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, pages 69–85. Springer, 2016.

Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021.

Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I 13*, pages 818–833. Springer, 2014.

Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. SWAG: A large-scale adversarial dataset for grounded commonsense inference. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 93–104, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1009. URL `https://aclanthology.org/D18-1009`.

Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6720–6731, 2019.

Chi Zhang, Feng Gao, Baoxiong Jia, Yixin Zhu, and Song-Chun Zhu. Raven: A dataset for relational and analogical visual reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5317–5327, 2019a.

Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. Vinvl: Revisiting visual representations in vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5579–5588, 2021.

Yundong Zhang, Juan Carlos Niebles, and Alvaro Soto. Interpretable visual question answering by visual grounding from attention supervision mining. In *2019 ieee winter conference on applications of computer vision (wacv)*, pages 349–357. IEEE, 2019b.

Jianxing Zheng, Zifeng Qin, Suge Wang, and Deyu Li. Attention-based explainable friend link prediction with heterogeneous context information. *Information Sciences*, 597:211–229, 2022.

Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning deep features for scene recognition using places database. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. URL `https://proceedings.neurips.cc/paper/2014/file/3fe94a002317b5f9259f82690aeea4cd-Paper.pdf`.

Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017.

Wanzheng Zhu and Suma Bhat. GRUEN for evaluating linguistic quality of generated text. *arXiv*, 2010.02498, 2020. ISSN 23318422. doi: 10.18653/v1/2020.findings-emnlp.9.

Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. Texygen: A benchmarking platform for text generation models. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 1097–1100, 2018.

Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7w: Grounded question answering in images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4995–5004, 2016.