

# How and where does CLIP process negation?

**Vincent Quantmeyer\***  
Utrecht University  
v.quantmeyer@gmail.com

**Pablo Mosteiro**  
Utrecht University  
p.mosteiro@uu.nl

**Albert Gatt**  
Utrecht University  
a.gatt@uu.nl

## Abstract

Various benchmarks have been proposed to test linguistic understanding in pre-trained vision & language (VL) models. Here we build on the existence task from the VALSE benchmark (Parcalabescu et al., 2022) which we use to test models’ understanding of negation, a particularly interesting issue for multimodal models. However, while such VL benchmarks are useful for measuring model performance, they do not reveal anything about the internal processes through which these models arrive at their outputs in such visio-linguistic tasks. We take inspiration from the growing literature on model interpretability to explain the behaviour of VL models on the understanding of negation. Specifically, we approach these questions through an in-depth analysis of the text encoder in CLIP (Radford et al., 2021), a highly influential VL model. We localise parts of the encoder that process negation and analyse the role of attention heads in this task. Our contributions are threefold. We demonstrate how methods from the language model interpretability literature (such as causal tracing) can be translated to multimodal models and tasks; we provide concrete insights into how CLIP processes negation on the VALSE existence task; and we highlight inherent limitations in the VALSE dataset as a benchmark for linguistic understanding.

## 1 Introduction

Research in vision & language (VL) modelling has produced various pre-trained models that are capable of jointly processing image and text information by learning multimodal representations (e.g.,

Li et al., 2019; Lu et al., 2019; Radford et al., 2021; Jia et al., 2021; Li et al., 2021). This makes them applicable to a host of downstream tasks, such as visual question answering, image caption generation or zero-shot image classification.

Various benchmarks have been proposed to test these models’ understanding of different linguistic features, such as word order (Akula et al., 2020), verb meaning (Hendricks and Nematzadeh, 2021), and compositionality (Thrush et al., 2022). The VALSE benchmark (Parcalabescu et al., 2022) was introduced to test these models’ ability to ground features such as existence, plurality, or spatial relations in images. An example of the existence piece is shown in Figure 1. Given an image, a model must choose between a correct caption and an incorrect foil, one of which contains a negation operator.

As such, this piece can be used to test a model’s understanding of negation, a particularly interesting issue for multimodal models, which typically include a visual backbone pretrained on computer vision tasks such as object labelling. The models themselves are further pretrained on image-text pairs where there is likely to be a *positive* bias, since captions describing images will typically refer to what is depicted there. This raises the question whether VL models are capable of processing operators such as “no” in instances such as those in Figure 1. Indeed, negation remains a weakness of even the most state-of-the-art large language models (Truong et al., 2023)

In line with these intuitions, initial VALSE results reveal that models only achieve moderate performance in this (and other) linguistic categories.

\*Work carried out as M.Sc. student at Utrecht University

However, while VL benchmarks such as VALSE are useful for measuring current and future model performance, they do not reveal anything about the internal processes through which these models arrive at their outputs in such visio-linguistic tasks.

We aim to make use of the growing literature on model interpretability (Räuker et al., 2022) in order to explain the behaviour (and shortcomings) of VL models on the understanding of negation. To do this, we use the existence sub-task in VALSE, with some extensions, exploiting localisation techniques to quantify the roles that different model components play in this task. This yields the following research question: *Which components of VL models are responsible for the model’s understanding of negation?* We address two issues that arise from this general question, namely (1) the extent to which processing of negation is localised vs. distributed; and (2) whether model performance on VALSE-like tasks involving negation can in part be explained by high-level dataset features.

Specifically, we approach these questions through an in-depth analysis of CLIP (Radford et al., 2021), a highly influential VL model. CLIP has a relatively simple design based exclusively on Transformers, which allows us to leverage interpretability techniques that target this architecture. Additionally, prior work by Parcalabescu and Frank (2023) shows that CLIP makes balanced use of text and image input and avoids so-called unimodal collapse (Madhyastha et al., 2018; Hessel and Lee, 2020; Frank et al., 2021), an important consideration for a study of multimodal model interpretability. Finally, CLIP remains central to developments in both vision (e.g. the CLIPSeg segmentation model; Lüddecke and Ecker, 2022) and VL tasks (e.g. CLIP is a component of several text-to-image and image-to-text models, including Mokady et al., 2021; Li et al., 2023; Ramesh et al., 2022; Rombach et al., 2022, among others).

In our analysis of negation, we focus on the CLIP text encoder. However, it is important to note that CLIP is pretrained with a multimodal contrastive objective, which has been shown to yield different representations compared to text-only encoders with comparable architecture but different training objectives (Wolfe and Caliskan, 2022). Thus, we take the insights into the text encoder’s ability to process negation as reflecting on the success or otherwise of the contrastive, multimodal pretraining in such models.

The contributions of this work are threefold:

firstly, we demonstrate how methods from the language model interpretability literature (e.g., causal tracing; Meng et al., 2023) can be translated to multimodal models and tasks; secondly, we provide concrete insights into how CLIP processes negation on the VALSE existence task; thirdly, we highlight inherent limitations in the VALSE dataset as a benchmark for linguistic understanding.

## 2 Related work

**Vision-and-language models** VL pretraining gained impetus from the development of multimodal, pretrained encoders inspired by BERT (Devlin et al., 2019). Bugliarello et al. (2021) provide a unified analysis of the varying VL BERT architectures.

With the introduction of CLIP (Radford et al., 2021), contrastive learning objectives have become prominent in VL models, with or without additional objectives that address multimodal fusion (Jia et al., 2021; Li et al., 2021; Singh et al., 2022a; Zeng et al., 2022). Models such as BLIP (Li et al., 2022) and FLAVA (Singh et al., 2022b) combine contrastive objectives with unimodal pretraining of vision and language encoders. Architectures such as Flamingo (Alayrac et al., 2022) and BLIP-2 (Li et al., 2023) reduce training cost by training relatively small networks to map between representations from pretrained image and language models.

In CLIP, an image encoder and a text encoder process their respective inputs completely separately from each other, i.e., without any multimodal cross-attention and project them into the same latent space. The goal of contrastive learning is to maximise similarity between matching image-text pairs, minimising the similarity between non-matching pairs. During inference, CLIP computes the similarity of an image and a text in the form of a scaled dot product between their embeddings. Contrastive objectives have been shown to yield better embedding representations (Wolfe and Caliskan, 2022) leading to improved performance on semantic evaluation tasks (Mu et al., 2018).

**Vision-and-language benchmarks** VL benchmarks focusing on specific linguistic phenomena play an important role in highlighting strengths and weaknesses in models’ grounding capabilities. For example, a recent study combining several benchmarks (Bugliarello et al., 2023) showed that models still find certain linguistic phenomena challenging, and that grounding capabilities may be less related



Type	Image	Caption	Foil
Negation in foil		There are giraffes	There are no giraffes
Negation in caption		There are no people	There are people

Figure 1: Examples from VALSE existence (Parcalabescu et al., 2022). Caption and foil only differ in the presence or absence of the negator “no”. The negator is either in the caption or the foil.

to model size, and more to other variables, including the fine-grained object recognition capabilities of the visual backbone (e.g. Zheng et al., 2022).

One class of benchmarks focuses on the robustness of models to syntactic permutations and/or their ability to reason compositionally when predicting whether visual inputs correspond to linguistic descriptions (e.g., Akula et al., 2020; Hendricks and Nematzadeh, 2021; Thrush et al., 2022; Ma et al., 2023; Yuksekgonul et al., 2023; Chen et al., 2023). Some of these benchmarks focus on specific linguistic phenomena, such as spatial relations (Liu et al., 2023; Kamath et al., 2023) or temporal relations (e.g. Kesen et al., 2024).

VALSE (Parcalabescu et al., 2022), on which we build the present study, prompts a model with an image along with both its correct caption and a foiled caption and tests a model’s ability to distinguish the caption from foil. This extends the original foiling task introduced by Shekhar et al. (2017). VALSE is divided into six sub-tasks or ‘pieces’, corresponding to six different linguistic phenomena. In this paper, we focus exclusively on the existence piece; see Figure 1.

**Model interpretability** R auker et al. (2022) define inner interpretability methods as those that help understand a model’s internal structures and activations. One recurring strategy in such techniques is to analyse the effect of perturbations or ablations on the model’s behaviour and output, whether this is applied to individual neurons (e.g. Zhou et al., 2018; Ghorbani and Zou, 2020) or to weights, with the goal of identifying modular subnetworks (Csord as et al., 2021).

The choice of a suitable level of granularity at which to apply ablation is largely dictated by the model’s size and complexity. Interpretability methods for transformers often operate at the level of attention heads, MHA modules, MLPs, or full Transformer layers.<sup>1</sup> Meng et al. (2023) introduced the causal tracing methodology to localise factual associations in a model.

In Meng et al. (2023), this localisation step serves as the basis for subsequent model editing in the ROME method. However, follow-up work has suggested that the ability to edit knowledge in a particular layer does not imply that this knowledge is localised in this layer (Hase et al., 2023) and can also introduce unwanted side effects (Hoelscher-Obermaier et al., 2023). Given these uncertainties surrounding model editing techniques, the present study focuses on localisation only.

A final line of relevant interpretability literature focuses on attention patterns in large Transformer models, which reveal the role of specific attention heads in processing linguistic phenomena such as syntactic roles (Clark et al., 2019; Kovaleva et al., 2019; Vig and Belinkov, 2019). All of these studies converge on the finding that pre-trained Transformer language models allocate significant attention to tokens that do not carry inherent semantic meaning, such as the separator token in BERT or the start-of-sequence token in GPT-2.

<sup>1</sup>Goh et al. (2021) also produced neuron-level interpretations of CLIP’s image encoder, albeit the ResNet and not the ViT variant.

	Correct	Ambiguous	Incorrect
	$d > 1$	$1 \geq d > -1$	$d \leq -1$
<b>Caption</b>	72	150	28
<b>Foil</b>	81	145	14

Table 1: Number of instances per segment in the VALSE existence dataset.

### 3 Methods

#### 3.1 Definitions

A forward pass in CLIP of a single VALSE existence instance (Fig. 1) consists of a text caption, a text foil, and an image. This produces one similarity score for caption and image and one for caption and foil, denoted  $S_{c,i}$  and  $S_{f,i}$ , respectively.

CLIP is said to correctly classify a caption-foil-image triple if  $S_{c,i} > S_{f,i}$ . We can quantify CLIP’s classification performance using the difference between the two similarities. We denote this classification score  $d = S_{c,i} - S_{f,i}$  and the absolute size of  $d$  can be seen as an indicator of CLIP’s confidence in the classification.

#### 3.2 Data

The VALSE existence benchmark consists of 505 image-caption-foil triples. The dataset is divided into instances where the negation is in the foil (249) and instances where the negation is in the caption (256). The presence or absence of a negation operator means that sometimes captions or foils can differ in token length. For our purposes, it is important that strings are of equal length; hence we insert the word *some* before the noun in non-negated sentences. See Appendix A.1 for full details.

CLIP only achieves a moderate accuracy of 0.686 on VALSE existence. To identify patterns of processing in the model that give rise to correct classification of negation it is necessary to analyse correctly and incorrectly classified instances separately. To do this consistently across different analyses, the dataset was divided into three segments (correct, ambiguous, incorrect) based on the classification score  $d$ . Table 1 shows the distribution of instances per segment.

#### 3.3 Causal tracing

Here we outline our adaption of the causal tracing method from Meng et al. (2023) for the part of the dataset where the negation is in the foil. Figure 2 provides a visual summary of the method.

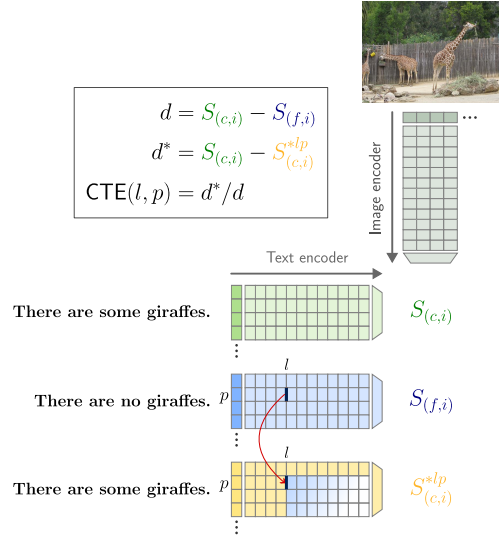


Figure 2: Illustration of the causal tracing methodology. The activation at a single position and layer from the negated forward pass are inserted into the corresponding layer and position of the non-negated forward pass. This shows what proportion of the original effect can be restored by this layer-position pair. Image and text are taken from VALSE existence (Parcalabescu et al., 2022).

A standard forward pass is carried out with caption, foil, and image, yielding the regular classification score  $d = S_{c,i} - S_{f,i}$ . Importantly, the activations from the forward pass at each layer and each position in the text encoder are recorded. In the subsequent modified forward pass only the (non-negated) caption is used in the forward pass alongside the image. During this forward pass, the text encoder’s activation at a given layer and position is replaced by the activation from the foil’s original forward pass at the corresponding layer and position. This is done individually for each combination of layer and position.

Intuitively, this achieves the following. The model processes the non-negated caption, but at a given layer and position it is made to behave as if it was processing the negated foil. If, and only if, a certain layer and position is specialised in processing negation, then substituting the activation from the negated forward pass into the non-negated one should affect the output in a visible way.

This intuition is quantified in the following way. For a given layer  $l$  and a position  $p$  the modified forward pass produces a similarity score  $S_{c,i}^{*lp}$ . This allows us to calculate a modified classification score

$$d^* = S_{c,i} - S_{c,i}^{*lp}$$

With this modified classification score we calculate

the causal tracing effect of layer  $l$  at position  $p$

$$\text{CTE}(l, p) = d^* / d$$

This effect represents the proportion of the original classification score  $d$  that can be “restored” by layer  $l$  at position  $p$ .

To apply this method to cases where the negation is in the caption, one has to swap caption and foil such that, once again, the activations from the negated sentence (now the caption) are substituted into the forward pass of the non-negated sentence (now the foil). This means that we obtain a modified classification score, which is used to calculate the causal tracing effect in the same way.

$$d^* = S_{f,i}^{*lp} - S_{f,i}$$

This method yields a causal tracing effect for each layer and position for each VALSE existence triple. All captions in the dataset share the same beginning (SOT, There, is/are, a/some, subject) and ending set of tokens (., EOT). However, they differ in the number of tokens in between these two sets. Therefore, the CTE from all positions in between the beginning and end sets of tokens are averaged into one placeholder position called “further subject tokens”. If there are no positions between the beginning and end sets, then a CTE of 0 is recorded at this position. Consequently, we can average CTEs across the dataset (or a segment thereof). To represent each instance according to its sequence length, the averaged effect at the “further subject tokens” position is weighted by the number of tokens that make up this position in each instance.

Lastly, we want to be able to describe the degree of localisation in particular layers. Localisation is strongest when one position in a layer, to the exclusion of all other positions, restores the full effect. Conversely, localisation is absent when each position restores the same proportion of the effect. Hence, we can quantify the degree of localisation in a layer  $l$  as the standard deviation of the causal tracing effects at each position in this layer, starting at the negator position.

### 3.4 Negator-selective attention in text encoder

The purpose of this analysis is to identify attention heads in CLIP’s text encoder that selectively pay attention to negators. Since a regular forward pass consists of both caption and foil, this yields two attention maps per head in the text encoder. Each attention map is an array of size  $P \times P$  where  $P$

is the number of positions in the input sequence, where the attention mask forces all elements to the right of the diagonal of this array to be 0.

The attention map is filtered to the column representing the position of the negator in the negated input sentence (or the quantifier in the corresponding non-negated sentence). To identify negator-selective attention, we subtract the values from the non-negated sentence from those from the negated sentence. Finally, the maximum of the resulting difference values is taken over all source positions and this represents the amount of negator-selective attention of a particular attention head on this particular dataset instance. This procedure can then be repeated over the whole dataset yielding an average negator-selective attention value  $a_{lh}^N$  for each attention head  $h$  in each layer  $l$ .

Instead of taking the maximum value over source positions, negator-selective attention can also be calculated for each source position. In heads with high negator-selective attention, this creates a more fine-grained picture of the negator-selective attention patterns involved.

To test the validity of the results from this analysis, it is further adapted to a subset of the CANNOT dataset (Anschütz et al., 2023), from which we use 554 negated sentences and create a positive counterpart for each. See Appendix A.2 for details.

## 4 Results

### 4.1 Causal tracing in text encoder

The left heatmap in Figure 3 shows the causal tracing effect per layer and position for the correct segment of the data with negation in the *foil*.

We are interested in the effect of components that lie in between the negator position in layer 0 (embeddings) and the last position in the final layer (encoder output), as these possibly mediate CLIP’s correct processing of negation in the text input.<sup>2</sup> Figure 3 shows that this effect is limited to only a subset of positions and layers and seems to suggest a path through the model. In particular, in layers 0-3 the effect is practically limited to the negator position, suggesting that in these early layers the negation information is processed mainly at its original position. The effect at the negator position then drops sharply at layer 4 and further decreases until the final layer. This indicates that

<sup>2</sup>Since the encoder uses masked attention, positions prior to the negator position cannot be affected by the intervention and therefore do not show any effect.

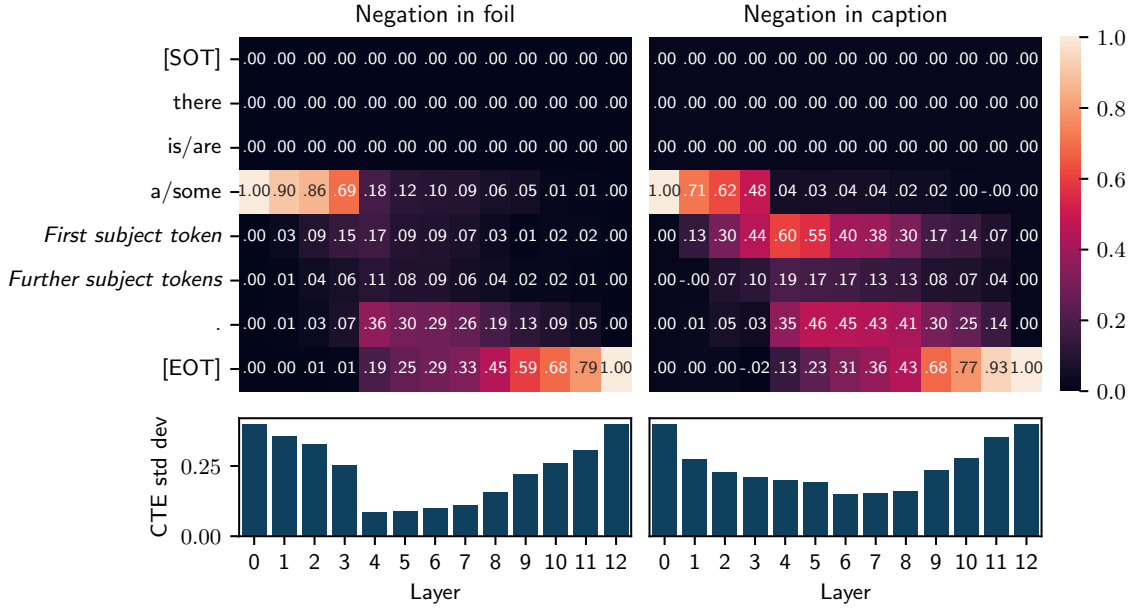


Figure 3: Causal tracing effect (CTE) of the correct segment, split by whether negation is in foil or caption. The heatmaps show the CTE of each layer-position pair in the text encoder. The bar charts show the standard deviation of all CTE in the corresponding layer as an overall measure of localisation. Layer 0 denotes the embedding layer.

the negator position only plays a pivotal role in the early layers and that the processing is in fact shifted to the second-to-last and last positions at layer 4. We will return to this in the analysis of attention patterns in Section 4.2. In the central layers 4-7 these two positions seem to play an equally important role, judging by their respective CTE, and from layer 8 onwards, the effect is concentrated in the last layer.

The bar charts in Figure 3 show the degree of localisation in each layer as measured by the standard deviation of the CTE. In line with the interpretation above, localisation is high in the early layers 0-3, then drops sharply in layer 4, remains low in the middle layers, and goes up again in the late layers 9-12.

The right part of Figure 3 shows the results from the same experiment on the correct segment of the data where the negation is in the *caption*. The general pattern of these results is comparable to the one described above. However, the first subject position already has a visible effect in the early layers, leading to reduced localisation. The effect of the first subject position becomes most pronounced in the middle layers which constitutes the most substantial difference between the two sets of results and in fact leads to greater localisation in the middle layers. In the late layers 9-12, the effect is once again concentrated in the last position.

## 4.2 Negator-selective attention in text encoder

Figure 4 shows the negator-selective attention of each attention head of each layer in CLIP’s text encoder, divided by whether the negation is in foil or caption. As expected, the patterns in both parts of the dataset are practically identical, since this analysis is not affected by any visual input. As a general observation, only a small subset of heads display any negator-selective attention (8% of heads with  $a_{ih}^N > 0.1$ ) and the majority of them are found in the early layers. The most negator-selective attention head is found in layer 4.

Note that these results are reported across all dataset segments (incorrect, ambiguous, correct), since the patterns do not meaningfully differ between them. This suggests that negator-selective attention cannot explain the difference in CLIP’s classification performance on different instances of VALSE existence, since the same patterns are found in correctly and incorrectly classified cases. In fact, none of the attention heads that show negator-selective attention of at least 0.1 show a correlation between negator-selective attention and classification score (all  $|r| < 0.2$ ).

Layer 4, where the most negator-selective attention is found, is the same layer, where the causal tracing results from Section 4.1 suggested that negation information is moved from its original position to later positions, in particular the second-to-last

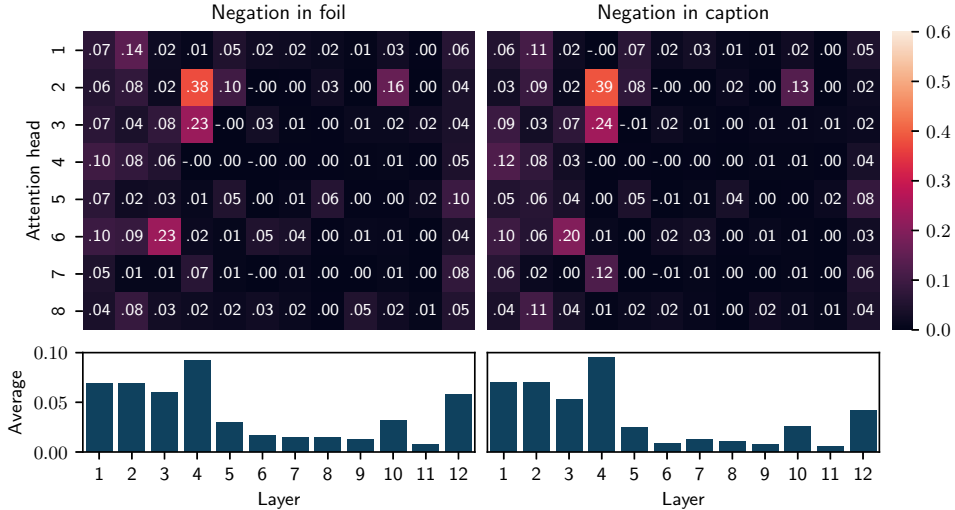


Figure 4: Negator-selective attention across all dataset segments, split by whether negation is in foil or caption. The heatmaps indicate the degree of negator-selective attention for each attention head and layer. The bar charts show the average of each layer as an overall measure of negator-selective attention.

one. We analyse the source of this negator-specific attention, i.e., which specific positions attend particularly to the negator position in the identified heads of interest. Figure 6 (Appendix A.3) confirms that the source of negator-selective attention in Head 2 is the second-to-last position. Furthermore, when the negation is in the caption, we find that additional negator-selective attention comes from the first subject position, which aligns with the greater role this position plays in this part of the dataset, as already suggested by the causal tracing results in Section 4.1. Thus, the causal tracing and negator-selective attention results form a coherent narrative.

We validate these observations using the CAN-NOT dataset. Here, we observe similar trends, with most negator-selective attention found in the early layers 1-4. See Appendix A.4 for details.

### 4.3 Dataset features

We investigate whether the similarity between a caption and a foil for a given VALSE instance is correlated with the instance’s classification score. Full details are in Appendix A.5, especially Figure 8. We make two primary observations. First the classification score is weakly correlated with the similarity between caption and foil, especially for those instances when the negation is in the foil. Second, longer sequences exhibit greater foil-caption similarity, leading to lower scores.

To investigate the effect of the size of the caption’s subject (e.g. ‘giraffe’ in Figure 1), we find

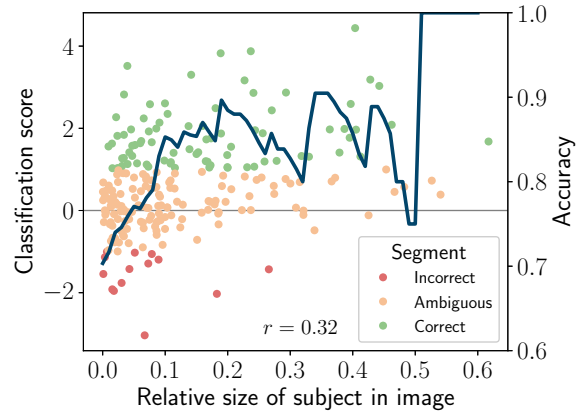


Figure 5: Relative size of image subject vs. CLIP’s classification score. All instances where the subject from the caption is shown in the image. Colour indicates dataset segment. The blue line shows classification accuracy when imposing a minimum subject size threshold.

its bounding box using CLIPSeg (Lüddecke and Ecker, 2022), and compare its relative size to the instance’s classification score (Figure 5). The correlation of  $r = 0.32$  shows that images with more prominent subjects tend to be classified more accurately. In fact, when imposing a subject size threshold of 0.1 (which removes 43% of instances), CLIP achieves an accuracy of 0.85. The accuracy as a function of the subject size threshold is shown by the line in Figure 5. Note, however, that the validity of these results decreases with higher thresholds, as the remaining sample size gets very small. Nonetheless, these results suggest that CLIP ex-

hibits better existence classification results on instances with more salient subjects.

## 5 Discussion

The causal tracing results from Section 4.1 suggest relatively strong localisation in the early (1-3) and late (8-12) layers, meaning that negation is largely represented at singular positions in these layers.

In layer 4, the CTE at the negator position drops sharply, and this coincides with the finding of negator-selective attention heads in layer 4 which appear to shift negation information to later positions. The locations of these attention heads also overlap with those found on the CANNOT dataset, which provides initial evidence that the CLIP text encoder uses certain attention heads for specific syntactic functions.

In the middle layers localisation is generally lower with no single position restoring more than 60% of the original effect. This implies that representation of negation is distributed across positions and that the model relies on *combining* the representations at each position in order to make correct judgements about negations.

Furthermore, the first subject token position appears to play a unique role in cases with negation in the caption, which could be due to the asymmetry in the two tasks. When the negation is in the foil, the label’s subject is shown in the image and, intuitively, once it is detected, a decision can be made and no further processing is necessary. Conversely, when the negation is in the caption, the entire image needs to be scanned to ensure that the label’s caption is in fact absent from all parts of the image. This difference could be part of the reason why the first subject token position appears to play a role up until deeper layers of the network, when the negation is in the caption. The effects of the subject position in deeper layers could imply that the subject information is in fact more deeply processed and thus more strongly represented in the final text encoder’s output which, in turn, could be conducive to the model’s task of “searching” for the subject in the image’s representation. However, these explanations are speculative and must not be accepted without further experiments.

Section 4.3 highlights that the label’s length and the subject’s size in the image show non-negligible correlations with respect to the classification score. This suggests that CLIP is better at the VALSE Existence task when labels are shorter and therefore

produce less similar multimodal embeddings and when the subject in the image is sufficiently large.

Arguably, the more variance in classification score can be explained on the basis of such dataset variables, the less CLIP’s benchmark score can be interpreted as an indicator of its linguistic understanding, thus calling into question the validity of the VALSE benchmark. However, none of the correlations found in the present study are particularly high and thus further analyses are needed to support this conclusion.

## 6 Limitations and future work

The degree of localisation found in CLIP’s text encoder is hard to interpret without reference to other results. Future work could extend the present methodology to other tasks and potentially other models.

Our study is also limited to simple effects of individual layer/position pairs. An analysis of the *interaction* of certain layers or positions (e.g., by simultaneously patching activations in multiple places during causal tracing) might draw a more robust and conclusive picture of the inner processes that govern CLIP’s understanding of negation.

More generally, localisation methods may not be suited for analysing model behaviour that is shown with only moderate reliability. Note that the methods used in the present study had originally been proposed and applied to language model capabilities that are shown reliably across a large corpus of data, e.g., indirect object identification (Wang et al., 2022), simple factual knowledge (Meng et al., 2023), or docstring completion (Heimersheim and Janiak, 2023). By contrast, CLIP does not reliably handle negation in a multimodal context (CLIP’s accuracy is only 66.9%) and these results are based on a relatively small dataset ( $n = 490$ ). In this case, methods like causal tracing do not intuitively lend themselves to *comparing* situations evincing a particular model behaviour to those where the behaviour is absent. That is because they focus on the degree to which an effect that represents a particular model behaviour can be restored or ablated, but the methodology breaks down when this effect isn’t present in the first place.

Thus, whilst illuminating the role of various components in CLIP’s processing of negation, we cannot provide strong insights into why this processing yields correct classifications only in a fraction of cases. Furthermore, since correct classification



only occurs in a subset of instances of VALSE, which is moderately sized to begin with, the results described here require a larger and potentially more diverse dataset to obtain greater validity.

With respect to the validity of the underlying VALSE benchmark, it might be worthwhile to conduct a larger study on dataset features (e.g., image brightness, contrast, etc.) that correlate with benchmark performance. Comparisons with other VL benchmarks would further help putting these results into perspective. Such features that are predictive of benchmark performance limit the validity of linguistic benchmarks and highlight variables that should be controlled for in the creation of future benchmarks.

## References

- Arjun Akula, Spandana Gella, Yaser Al-Onaizan, Song-Chun Zhu, and Siva Reddy. 2020. [Words Aren't Enough, Their Order Matters: On the Robustness of Grounding Visual Referring Expressions](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6555–6565, Online. Association for Computational Linguistics.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. 2022. [Flamingo: A Visual Language Model for Few-Shot Learning](#). *Preprint*, arxiv:2204.14198.
- Miriam Anshütz, Diego Miguel Lozano, and Georg Groh. 2023. [This is not correct! Negation-aware Evaluation of Language Generation Systems](#). *Preprint*, arxiv:2307.13989.
- Emanuele Bugliarello, Ryan Cotterell, Naoaki Okazaki, and Desmond Elliott. 2021. [Multimodal Pretraining Unmasked: A Meta-Analysis and a Unified Framework of Vision-and-Language BERTs](#). *Preprint*, arxiv:2011.15124.
- Emanuele Bugliarello, Laurent Sartran, Aishwarya Agrawal, Lisa Anne Hendricks, and Aida Nematzadeh. 2023. [Measuring Progress in Fine-grained Vision-and-Language Understanding](#). *Preprint*, arxiv:2305.07558.
- Xinyi Chen, Raquel Fernández, and Sandro Pezzelle. 2023. [The BLA Benchmark: Investigating Basic Language Abilities of Pre-Trained Multimodal Models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5817–5830, Singapore. Association for Computational Linguistics.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [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. Association for Computational Linguistics.
- Róbert Csordás, Sjoerd van Steenkiste, and Jürgen Schmidhuber. 2021. [Are Neural Nets Modular? Inspecting Functional Modularity Through Differentiable Weight Masks](#). *Preprint*, arxiv:2010.02066.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#). *Preprint*, arxiv:1810.04805.
- Stella Frank, Emanuele Bugliarello, and Desmond Elliott. 2021. [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 (EMNLP'21)*. ArXiv: 2109.04448.
- Amirata Ghorbani and James Zou. 2020. [Neuron Shapley: Discovering the Responsible Neurons](#). *Preprint*, arxiv:2002.09815.
- Gabriel Goh, Nick Cammarata †, Chelsea Voss †, Shan Carter, Michael Petrov, Ludwig Schubert, Alec Radford, and Chris Olah. 2021. [Multimodal Neurons in Artificial Neural Networks](#). *Distill*, 6(3):e30.
- Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. 2023. [Does Localization Inform Editing? Surprising Differences in Causality-Based Localization vs. Knowledge Editing in Language Models](#). *Preprint*, arxiv:2301.04213.
- Stefan Heimersheim and Jett Janiak. 2023. [A circuit for Python docstrings in a 4-layer attention-only transformer](#). *Alignment Forum*.
- Lisa Anne Hendricks and Aida Nematzadeh. 2021. [Probing Image-Language Transformers for Verb Understanding](#). *Preprint*, arxiv:2106.09141.
- Jack Hessel and Lillian Lee. 2020. [Does my multimodal model learn cross-modal interactions? It's harder to tell than you might think!](#) In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 861–877, Online. Association for Computational Linguistics.
- Jason Hoelscher-Obermaier, Julia Persson, Esben Kran, Ioannis Konstas, and Fazl Barez. 2023. [Detecting Edit Failures In Large Language Models: An Improved Specificity Benchmark](#). *Preprint*, arxiv:2305.17553.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. [spaCy: Industrial-strength Natural Language Processing in Python](#).

- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. 2021. [Scaling Up Visual and Vision-Language Representation Learning With Noisy Text Supervision](#). *Preprint*, arxiv:2102.05918.
- Amita Kamath, Jack Hessel, and Kai-Wei Chang. 2023. [What’s “up” with vision-language models? Investigating their struggle with spatial reasoning](#). *arXiv preprint*. ArXiv:2310.19785 [cs].
- Ilker Kesen, Andrea Pedrotti, Mustafa Dogan, Michele Cafagna, Emre Can Acikgoz, Letitia Parcalabescu, Iacer Calixto, Anette Frank, Albert Gatt, Aykut Erdem, and Erkut Erdem. 2024. [Vilma: A zero-shot benchmark for linguistic and temporal grounding in video-language models](#). In *Proceedings of the 12th International Conference on Learning Representations (ICLR)*, Vienna, Austria.
- Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. 2019. [Revealing the Dark Secrets of BERT](#). 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 4365–4374, Hong Kong, China. Association for Computational Linguistics.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. [BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models](#). *Preprint*, arxiv:2301.12597.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. [BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation](#). *Preprint*, arxiv:2201.12086.
- Junnan Li, Ramprasaath R. Selvaraju, Akhilesh Deepak Gotmare, Shafiq Joty, Caiming Xiong, and Steven Hoi. 2021. [Align before Fuse: Vision and Language Representation Learning with Momentum Distillation](#). *Preprint*, arxiv:2107.07651.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. [VisualBERT: A Simple and Performant Baseline for Vision and Language](#). *Preprint*, arxiv:1908.03557.
- Fangyu Liu, Guy Emerson, and Nigel Collier. 2023. [Visual Spatial Reasoning](#). *Preprint*, arxiv:2205.00363.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. [ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks](#). *Preprint*, arxiv:1908.02265.
- Timo Lüddecke and Alexander S. Ecker. 2022. [Image Segmentation Using Text and Image Prompts](#). *Preprint*, arxiv:2112.10003.
- Zixian Ma, Jerry Hong, Mustafa Omer Gul, Mona Gandhi, Irena Gao, and Ranjay Krishna. 2023. [CREPE: Can Vision-Language Foundation Models Reason Compositionally?](#) *Preprint*, arxiv:2212.07796.
- Pranava Madhyastha, Josiah Wang, and Lucia Specia. 2018. [Defoiling Foiled Image Captions](#). In *Proceedings of 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT’18)*, pages 433–438. ArXiv: 1805.06549.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2023. [Locating and Editing Factual Associations in GPT](#). *Preprint*, arxiv:2202.05262.
- Ron Mokady, Amir Hertz, and Amit H. Bermano. 2021. [ClipCap: CLIP Prefix for Image Captioning](#). *Preprint*, arxiv:2111.09734.
- Jiaqi Mu, Suma Bhat, and Pramod Viswanath. 2018. [All-but-the-Top: Simple and Effective Postprocessing for Word Representations](#). *Preprint*, arxiv:1702.01417.
- Letitia Parcalabescu, Michele Cafagna, Lilitta Muradjan, Anette Frank, Iacer Calixto, and Albert Gatt. 2022. [VALSE: A Task-Independent Benchmark for Vision and Language Models Centered on Linguistic Phenomena](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8253–8280, Dublin, Ireland. Association for Computational Linguistics.
- Letitia Parcalabescu and Anette Frank. 2023. [MM-SHAP: A Performance-agnostic Metric for Measuring Multimodal Contributions in Vision and Language Models & Tasks](#). *Preprint*, arxiv:2212.08158.
- 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. 2021. [Learning Transferable Visual Models From Natural Language Supervision](#). *Preprint*, arxiv:2103.00020.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. [Hierarchical Text-Conditional Image Generation with CLIP Latents](#). *Preprint*, arxiv:2204.06125.
- Tilman Räuher, Anson Ho, Stephen Casper, and Dylan Hadfield-Menell. 2022. [Toward Transparent AI: A Survey on Interpreting the Inner Structures of Deep Neural Networks](#). *Preprint*, arxiv:2207.13243.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. [High-Resolution Image Synthesis with Latent Diffusion Models](#). *Preprint*, arxiv:2112.10752.
- Ravi Shekhar, Sandro Pezzelle, Yauhen Klimovich, Aurélie Herbelot, Moin Nabi, Enver Sangineto, and Raffaella Bernardi. 2017. [FOIL it! Find One mismatch between Image and Language caption](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 255–265, Vancouver, Canada. Association for Computational Linguistics.

- Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. 2022a. [FLAVA: A Foundational Language And Vision Alignment Model](#). *Preprint*, arxiv:2112.04482.
- Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. 2022b. [FLAVA: A Foundational Language And Vision Alignment Model](#). In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15617–15629. ISSN: 2575-7075.
- Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. 2022. [Winoground: Probing Vision and Language Models for Visio-Linguistic Compositionality](#). *Preprint*, arxiv:2204.03162.
- Thinh Hung Truong, Timothy Baldwin, Karin Verspoor, and Trevor Cohn. 2023. [Language models are not naysayers: An analysis of language models on negation benchmarks](#). *Preprint*, arxiv:2306.08189.
- Jesse Vig and Yonatan Belinkov. 2019. [Analyzing the Structure of Attention in a Transformer Language Model](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 63–76, Florence, Italy. Association for Computational Linguistics.
- Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2022. [Interpretability in the Wild: A Circuit for Indirect Object Identification in GPT-2 small](#). *Preprint*, arxiv:2211.00593.
- Robert Wolfe and Aylin Caliskan. 2022. [Contrastive Visual Semantic Pretraining Magnifies the Semantics of Natural Language Representations](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3050–3061, Dublin, Ireland. Association for Computational Linguistics.
- Mert Yuksekogonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. 2023. [When and why vision-language models behave like bags-of-words, and what to do about it?](#) *Preprint*, arxiv:2210.01936.
- Yan Zeng, Xinsong Zhang, and Hang Li. 2022. [Multi-Grained Vision Language Pre-Training: Aligning Texts with Visual Concepts](#). In *Proceedings of the 39th International Conference on Machine Learning*, pages 25994–26009. PMLR.
- Kaizhi Zheng, Xiaotong Chen, Odest Chadwicke Jenkins, and Xin Wang. 2022. [Vlmbench: A compositional benchmark for vision-and-language manipulation](#). *Advances in Neural Information Processing Systems*, 35:665–678.
- Bolei Zhou, Yiyou Sun, David Bau, and Antonio Torralba. 2018. [Revisiting the Importance of Individual Units in CNNs via Ablation](#). *Preprint*, arxiv:1806.02891.

## A Appendix

### A.1 Preprocessing of VALSE instances

Since caption and foil in the VALSE existence dataset differ only in the presence of the negator, they sometimes have a different number of tokens. Concretely, this is the case in “bare plural” sentences where there is no article or other qualifier in the non-negated sentence (e.g., “There are tennis players.” vs “There are *no* tennis players.”). Identifying differences in how CLIP processes negated vs. non-negated labels is a core facet of the present study and such comparisons are greatly facilitated if caption and foil have the same number of tokens. Therefore, labels were rephrased to achieve equal sequence length by inserting the qualifier “some” into the non-negated plural sentences right before the subject. For example, “There are tennis players” was rephrased to “There are *some* tennis players”. 15 instances (0.03%) from the original dataset have labels that do not follow the simple “There is/are no [subject] ...” structure and therefore aren’t amenable to the rephrasing rule described above. For reasons of simplicity, these were omitted from the rephrased dataset.

Importantly, rephrasing the dataset in this way only led to minor changes in CLIP’s classification accuracy on this dataset (0.691 before, 0.686 after rephrasing). All analyses are based on the rephrased dataset, unless denoted otherwise.

### A.2 CANNOT dataset

We use the CANNOT dataset to independently validate our analysis of negator-selective attention in the CLIP text encoder.

For the present purposes, the dataset is filtered to 554 negated sentences that contain the word “no” as the determiner of the sentence subject (e.g., “Medical organizations recommend *no* alcohol during pregnancy for this reason”), using a tokeniser from the spacy Python library (Honnibal et al., 2020).

For each of these sentences, a non-negated counterpart is then generated by replacing the word “no” with “some”.

This yields a set of sentence pairs, comparable to the caption-foil pairs from VALSE existence, which thus allows us to apply the same methodology for negator-selective attention.

### A.3 Negator-selective attention on VALSE

As discussed in Section 4.2, Figure 6 confirms that the source of negator-specific attention in Head 2 is the second-to-last position.

### A.4 Negator-selective attention results on CANNOT

For validation purposes, Figure 7 shows negator-selective attention on a subset of the CANNOT dataset. Just like on the VALSE dataset, most negator-selective attention is found in the early layers 1-4. Head 2 in layer 4 once again shows particularly high negator-selective attention, albeit not the highest, which here is found in head 1 in layer 2. In summary, this provides converging evidence for the negator-selective attention results found in VALSE existence.

### A.5 Dataset features

Figure 8 shows the cosine similarity of each instance’s caption and foil in CLIP’s multimodal embedding space against that instance’s classification score, split by whether the negation is in the caption or foil.

When the negation is in the foil, similarity and score are weakly correlated ( $r = -0.22$ ), whereas no correlation is found when the negation is in the caption ( $r = 0.03$ ). The latter is however influenced by the presence of a set of outliers, all with the same caption “There are no people.”. Removing them from this analysis yields a correlation of  $r = -0.20$ , comparable to the one found when the negation is in the foil.

Figure 8 also encodes sequence length, with longer sequences (darker colour) tending to exhibit greater caption-foil similarity. This is to be expected since caption and foil differ in exactly one position. If the total number of positions increases, then the relative size of this difference decreases, leading to greater similarity. These results suggest that CLIP’s failure to correctly classify some VALSE Existence instances might be partly due to instances with longer captions and foils that are more similar in their representation and therefore more difficult to tell apart. However, filtering the dataset to instances with shorter sequences does not meaningfully improve CLIP’s accuracy, suggesting that sequence length plays a minor role at best.

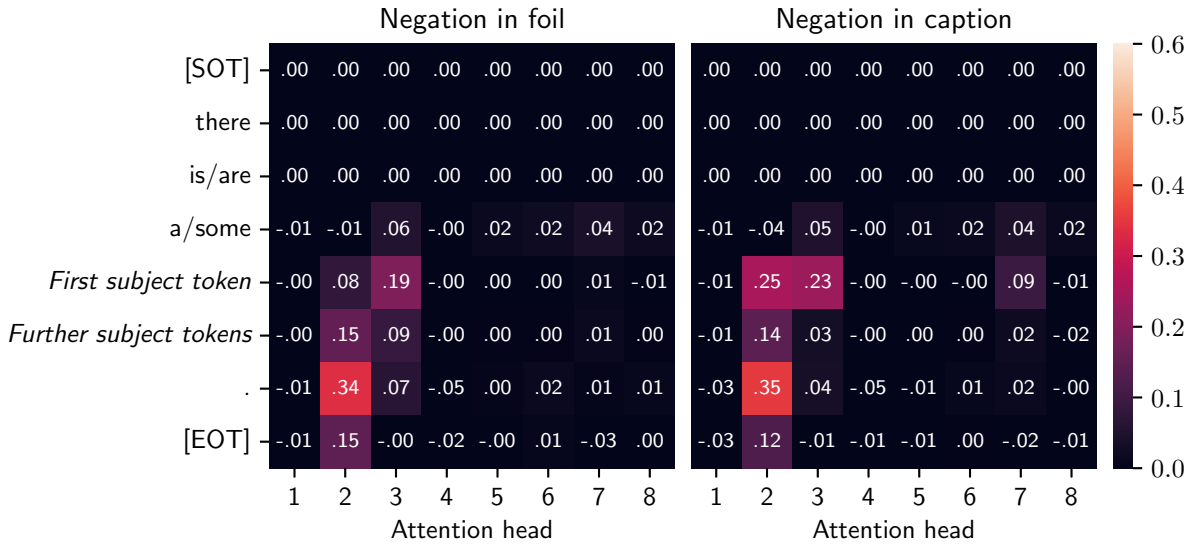


Figure 6: Source of negator-selective attention in layer 4 across all dataset segments, split by whether negation is in foil or caption. The heatmaps show the degree of negator-selective attention from each sequence position (y-axis) in each attention head (x-axis).

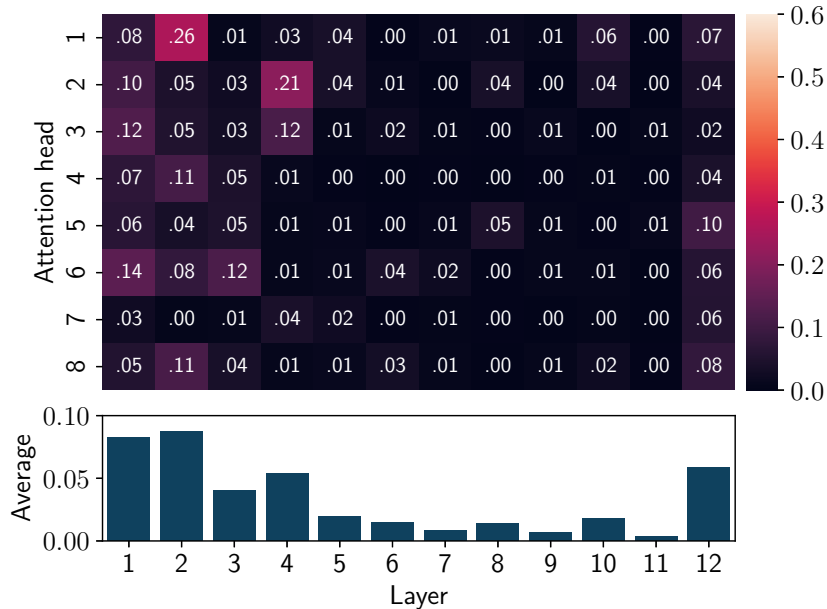


Figure 7: Negator-selective attention on the CANNOT dataset, to validate the results from Figure 4. The heatmaps indicate the degree of negator-selective attention for each attention head and layer. The bar charts show the average of each layer as an overall measure of negator-selective attention.

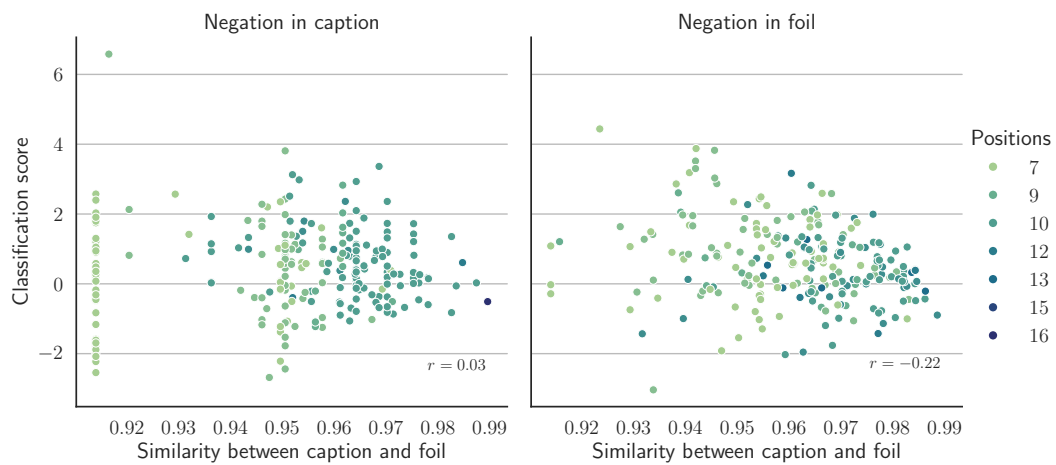


Figure 8: Cosine similarity of caption and foil in CLIP’s multimodal embedding space vs. CLIP’s classification score. Colour indicates dataset sequence length (i.e., number of tokens in sequence).