SPECS: Specificity-Enhanced CLIPScore for Long Image Caption Evaluation

Abstract

As interest grows in generating long, detailed image captions, standard evaluation metrics become increasingly unreliable. N-gram-based metrics though efficient, fail to capture semantic correctness. Representational Similarity (RS) metrics, designed to address this, initially saw limited use due to high computational costs, while today, despite advances in hardware, they remain unpopular due to low correlation to human judgments. Meanwhile, metrics based on large language models (LLMs) show strong correlation with human judgments, but remain too expensive for iterative use during model development. We introduce SPECS (Specificity-Enhanced CLIP-Score), a reference-free RS metric tailored to long image captioning. SPECS modifies CLIP with a new objective that emphasizes specificity: rewarding correct details and penalizing incorrect ones. We show that SPECS matches the performance of open-source LLMbased metrics in correlation to human judgments, while being far more efficient. This makes it a practical alternative for iterative checkpoint evaluation during image captioning model development. Our code can be found at https://github.com/mbzuai-nlp/SPECS.

1 Introduction

Image captioning has been a key topic in vision-language research, offering a controlled setting to study grounded language generation (Karpathy and Fei-Fei, 2015). While early efforts focused on short, general captions (Vinyals et al., 2015), recent work has shifted toward generating long, detailed descriptions that capture fine-grained visual information(Johnson et al., 2016; Cho et al., 2022; Doveh et al., 2023; Li et al., 2023). This complex task requires strong visual grounding and improved cross-modal alignment (Liu et al., 2024; Li et al., 2021; Xie et al., 2025). It expands the scope of generative vision-language modeling but also mag-

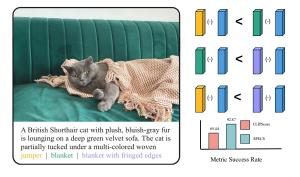


Figure 1: Example of specificity in caption evaluation. Given an image, captions with increasing correct details ("jumper" \rightarrow "blanket" \rightarrow "blanket with fringed edges") should receive progressively higher cosine similarity, (\cdot), to the image. SPECS ranks these minimal pairs correctly at a high rate, reflecting its strong specificity.

nifies a long-standing challenge: how to evaluate captions reliably and efficiently.

Automatic evaluation in captioning, as in other natural language generation tasks, has long been a challenge (Otani et al., 2023; Wang et al., 2023). Early metrics rely on n-gram overlap between generated captions and references. While computationally efficient, these methods fail to capture semantic similarity, often penalizing valid paraphrases and underestimating the severity of hallucinations, even in short captions (Papineni et al., 2002; Banerjee and Lavie, 2005; Vedantam et al., 2015; Lin, 2004).

To address these limitations, Representational Similarity (RS) metrics use pretrained vision-language models to compare image and caption embeddings in a shared feature space (Hessel et al., 2021; Sarto et al., 2023), often eliminating the need for reference captions. More recently, evaluation metrics based on large language models (LLMs) have become the standard, showing strong correlation with human judgments, especially as the length of generated image captions increases (Chan et al., 2023; Yu et al., 2024; Ye et al., 2025).

Each new generation of evaluation metrics

brings improvements in semantic expressiveness, but often at the cost of higher computational requirements. As a result, there is frequently a mismatch between what is technically feasible and what is practical for routine model development. For instance, although CLIPScore (Hessel et al., 2021) demonstrated stronger semantic alignment than traditional metrics like CIDEr (Vedantam et al., 2015), its adoption remained limited due to the high computational cost at the time of its release. A similar pattern is now emerging with LLM-based metrics: while they achieve state-of-the-art correlation with human judgments, their inference cost makes them impractical for iterative evaluation during model training and development.

This tradeoff between reliability and efficiency becomes particularly problematic in the context of long image captioning, where no existing metric can strike a reasonable balance between the two. Recent studies have shown that even simple metrics like BLEU-4 outperform RS metrics like CLIPScore in terms of correlation with human judgments (Ye et al., 2025, sample-level Kendall's Tau correlations of 0.27 and 0.17, respectively). These results underscore a pressing gap: there is no metric that reliably evaluates the quality of long captions while remaining computationally practical.

In this work, we introduce **SPECS** (**SPecificity-Enhanced CLIP-Score**), a reference-free RS metric designed for the evaluating long image captions. SPECS builds on a long-context adaptation of CLIP (Zhang et al., 2024) and incorporates a new training objective that emphasizes *specificity*: rewarding correct details and penalizing incorrect ones. In extensive evaluations, SPECS matches the best open-source LLM-based metric (Lee et al., 2024) in terms of correlation to human judgments, while being over two orders of magnitude more efficient. SPECS is a practical and scalable solution for iterative model development in long caption generation.

2 Related work

Image captioning metrics are central to evaluating vision language models (VLMs). These metrics score how well a model can describe an image, in a way that aligns with human judgment.

Image Caption Evaluation Methods Early work relied on *n*-gram matching, where metrics such as BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), ROUGE (Lin, 2004), and CIDEr (Vedantam et al., 2015) com-

pared generated captions against human references using surface lexical overlap. While easy to compute, these metrics often fail when captions use different but valid phrasings, leading to low correlation with human judgment. To overcome this, later approaches explored semantic parsing, such as SPICE (Anderson et al., 2016), which evaluates scene-graph structures. CLIPScore (Hessel et al., 2021) and PACScore (Sarto et al., 2023) build on pre-trained vision—language models such as CLIP (Radford et al., 2021), assessing captions through representational similarity between images and their captions.

More recently, large language models (LLMs) have been used as evaluators by prompting them to score alignment between captions and images: FaithScore (Jing et al., 2023) extracts atomic facts from captions using an MLLM and verifying each fact against the image to measure faithfulness. CLAIR (Chan et al., 2023) leverages LLMs in a zero-shot setting to score captions and explain their judgments, GPT4-Eval (Liu et al., 2023) is evaluated by prompting GPT-4 to judge multimodal responses, with performance measured by relative quality scores. RLAIF-V (Yu et al., 2024) evaluates captions by prompting open-source MLLMs to verify decomposed claims and score hallucination reduction. FLEUR (Lee et al., 2024) prompts an open-source MLLM to directly compare an image with a candidate caption in a reference-free setting, producing both a numerical score and a natural-language rationale to align evaluation with human judgments. CAPTURE (Dong et al., 2024) evaluates captions by extracting and matching finegrained visual details. DCScore (Ye et al., 2025) evaluates captions by decomposing them into primitive information units and measuring their accuracy and coverage using GPT-4. These LLM-based methods show strong performance but remain prohibitively expensive and mostly closed-source.

CLIP-based adaptations Metrics based on representational similarity carry promise as they operate on the level of semantics rather than surface, and they are less costly than LLMs. However, CLIP, the common choice for base model in such metrics, is not suitable for long caption evaluation for two key reasons: it lacks compositionality and it can only take up to 77 subword tokens. To address these limitations, a line of research has modified CLIP in various ways. NegCLIP (Yuksekgonul et al., 2022) is the first work to address composi-

tionality, enhancing robustness by penalizing captions that introduce wrong attributes and reducing the risk of rewarding hallucinated details. Later, LaCLIP (Fan et al., 2023) and TripletCLIP (Patel et al., 2024) further advance compositional evaluation, making metrics more sensitive to how attributes and objects are combined. DAC (Doveh et al., 2023) fine-tunes CLIP on enhanced captions to improve caption density and strengthen compositional reasoning. DCI (Urbanek et al., 2024) introduces a densely captioned dataset with long, region-aligned descriptions and trains CLIP for improved long-text understanding. To address the issue of context length, LongCLIP (Zhang et al., 2024) modifies CLIP to process longer textual inputs through interpolation of the positional embeddings. LongCLIP exhibits strong long-caption understanding in retrieval tasks. The models listed above were not designed as metrics as such, but we hypothesize that their intended strengths might benefit long caption evaluation.

Recent evaluation metrics have made progress, but they still face two main problems. Some metrics do not capture detailed information well, while others do so at a high computational cost. We propose a new metric which leverages representational similarity, with a focus on specificity—the ability of a vision-language model to consistently prefer more informative, visually grounded captions at varying caption lengths.

3 Specificity

Let us consider the three caption variants depicted in Figure 1. Describing the cat as tucked under a blanket is correct, and because the caption already contains other relevant details, a good evaluation metric should assign it a high score. If the caption further specified that the blanket has *fringed edges*, the score should increase slightly, reflecting the correct additional detail. On the other hand, if instead of a *blanket* the caption said that the cat was lying under a *jumper*, that should result in a slightly lower score—most of the details remain accurate, but this particular object mentioned is incorrect. This simple example illustrates the notion of *speci*ficity which we adapt from Xu et al. (2024) to mean: the ability of a text representation to encode every detail in a caption in a way that correctly reflects the relevance of this detail to a reference image. A metric based on a specificity-enhanced model would thus favor captions that include more relevant details and penalize those that omit important information or introduce hallucinated or erroneous content. Such a metric implicitly implements the notions of soft precision and recall.

3.1 Detail Units

To concretely evaluate specificity, we begin by introducing the key concept of a **detail unit**. The abstract notion behind a detail unit refers to any minimal bit of information in a caption, such as the presence of a *blanket*, the *fringed edges* of the blanket, etc. For operational purposes, however, we define a detail unit to mean a phrase which contributes at least one new visual detail (and possibly more), and fits syntactically and semantically within the preceding context. Under this definition, a blanket is a detail unit, and so is a blanket with fringed edges, but a blanket with is not, and neither is The cat in the middle of the caption in Figure 1, since it does not contribute new information.

Formally, we denote an image-caption pair as $\{i, c\}$, and decompose a caption as c = $\{d_1, d_2, \dots, d_m\}$ where each d_i is a detail unit. Every subsequence of detail units, built cumulatively from left to right, constitutes a valid caption: $c_1 = \{d_1\}, c_2 = \{d_1+d_2\}, ..., \{c = d_1+\cdots+d_m\},\$ each containing progressively more information. Given a ground-truth, high-quality caption, this ordered sequence should exhibit monotonically increasing representational similarity to its reference image, under a specificity-enhanced model. Conversely, if an incorrect detail unit is added at any point, this should be reflected in a dip in the similarity score. This decomposition provides a structured way to test and enhance model specificity to visual detail across any caption length.

Detail units that contain relevant information are referred to as **positive** (d_+) , while detail units that introduce content not grounded in the image are referred to as **negative** (d_{-}) . The expected behavior of a model with good specificity is then to assign higher similarity to the pair $\{i, c_j + d_+\}$ than to the pair $\{i, c_i\}$, and a lower similarity to the pair $\{i, c_i + d_-\}$ than to the pair $\{i, c_i\}$, where $j \in [1,...,m]$ and c_j is a partial caption for the image. Each triplet, $\{i, c_i, c_i + d_+\}$ and $\{i, c_j, c_j + d_-\}$ constitutes a minimal pair of captions grounded in an image, the former being positive and the latter negative. Defining specificity with reference to both positive and negative details ensures that a model does not learn to simply assign higher similarity to longer captions, but evaluates

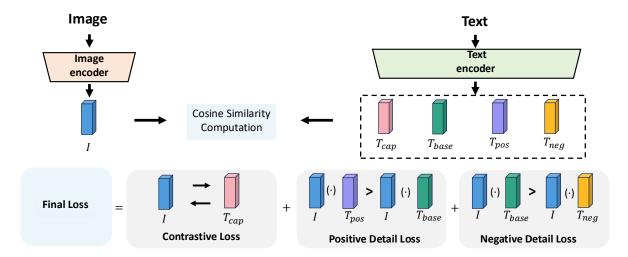


Figure 2: Training framework. Given an image and its caption, we produce a base caption, a more positive caption, and a negative caption. The model computes embeddings and is trained with three losses: a contrastive loss $\mathcal{L}_{\text{contrastive}}$ on the full caption, a positive detail loss \mathcal{L}_{pos} to prefer more informative descriptions, and a negative detail loss \mathcal{L}_{neg} to penalize misleading ones. This setup encourages sensitivity to fine-grained textual differences. The symbol (\cdot) denotes cosine similarity computation. Here, I is the image embedding, T_{cap} is the original caption, T_{base} is the base caption, T_{pos} is the more detailed caption, and T_{neg} is the negative caption.

the relevance of every new detail in the caption against the reference image.

3.2 Specificity Rate

To aggregate specificity across a set of minimal pairs, we introduce the **Specificity Rate** (SR). We define two variants: SR_{pos} measures the proportion of cases in which adding an additional relevant detail (positive detail unit) increases the similarity score with the image, while SR_{neg} measures the proportion of cases in which adding an irrelevant detail (negative detail unit) decreases the similarity. Given a set of N positive or negative triplets, we compute the SR as follows:

$$SR_{pos} = \frac{1}{N} \sum_{j}^{N} \mathbb{I}[\theta(i, c_j + d_+) > \theta(i, c_j)]$$
 (1)

$$SR_{neg} = \frac{1}{N} \sum_{j}^{N} \mathbb{I}[\theta(i, c_j) < \theta(i, c_j + d_-)]$$
 (2)

where $\mathbb{I}[\cdot]$ is the indicator function which outputs 1 if the condition inside is true and 0 otherwise, and θ is the cosine similarity between the representations of image and text. This formulation captures the rate at which representational similarity increases with added positive details, or decreases with added negative ones, thus measuring model specificity.

3.3 Specificity-Aware Learning

Although specificity can be used purely for evaluation, we can also enforce it during training. To encourage the model to prefer captions that describe images with greater relevant detail, we introduce a training objective that rewards higher similarity scores for incrementally more informative captions, and lower similarity scores for less accurate ones. Given a dataset of N positive and N negative triplets, we define the following hinge loss with a dynamic margin:

$$\mathcal{L}_{pos} = \frac{1}{N} \sum_{i}^{N} \max(0, \theta(i, c_j) - \theta(i, c_j + d_+) + \epsilon), \quad (3)$$

where ϵ is a batch-wise average similarity difference between positive and base captions, which is detached from gradient computation and clamped for numerical stability:

$$\epsilon = \operatorname{detach}\left(\frac{1}{N}\sum_{i}^{N}\left(\theta(i, c_{j} + d_{+}) - \theta(i, c_{j})\right)\right),$$

The negative loss, \mathcal{L}_{neg} is computed by analogy, from the negative triplets in the dataset.

The final training objective combines the contrastive loss, the positive detail loss and the negative detail loss:

$$\mathcal{L} = \alpha \mathcal{L}_{\text{contrastive}} + \beta \mathcal{L}_{\text{pos}} + \gamma \mathcal{L}_{\text{neg}}, \qquad (4)$$

where α , β , and γ are weighting hyperparameters tuned on a validation set. Figure 2 illustrates the overall training framework.

3.4 Metric Computation

Given a SPecificity-Enhanced CLIP (SPEC) model trained as described above, we score candidate captions, \hat{c} , against input image, i, as follows:

$$SPECS = \theta_{norm}(i, \hat{c})$$
 (5)

i.e., the metric uses cosine similarity, clipped at 0.

4 Experiments and Evaluation

4.1 Training and Validation Datasets

We train our model on the ShareGPT-4V dataset (Chen et al., 2024), which contains 1.2 million high-quality image-caption pairs synthetically generated by a strong captioning model, instructed to mention object attributes, spatial layouts, and aesthetic properties. The images in the dataset are sourced from COCO (Lin et al., 2014), SAM (Kirillov et al., 2023), and LAION (Schuhmann et al., 2022), and captions are 143 tokens long on average.

For intrinsic specificity evaluation, we use the sDCI dataset (Urbanek et al., 2024), consisting of 7805 images, each paired with 10 captions, which are synthetically designed to fit in CLIP's context window of 77 tokens. This underutilizes the full context window of our model, but enables controlled comparisons to other models, constrained by the 77-token context window. We compare against the models introduced in Section 2, as well as SigLIP (Zhai et al., 2023), a recent vision-language model, which replaces the softmax contrastive loss in CLIP with a pairwise sigmoid loss, improving efficiency, scalability and zero-shot transfer.

4.2 Data Preprocessing

To create the data needed to measure and train for specificity, we build a pipeline that segments captions into detail units.

Main Logic We considered various methods to segment captions into detail units, based on part-of-speech tagging, dependency parsing and coreference resolution: the results were either unsatisfactory, slow to obtain or obstructed by technical challenges with the deployment of outdated libraries. The solution that proved best in terms of speed, ease of implementation and quality was

obtained with the help of GPT-4. We presented the model with an example of a manually annotated caption and had it generate Python code that implements the segmentation pattern found in the example. The resulting code is based on part-of-speech tagging and a rule-based grammar (see Appendix A). Through manual inspection, we established that the solution is largely effective, but it somewhat oversegments the captions.

False Negatives vs. False Positives Considering the intended use of the segmented data, we determined that allowing for false negatives (i.e., missing splits) is less harmful than introducing false positives (i.e., incorrect extra splits). In other words, we prefer case (a), where a possible split is missing, over case (b), where an erroneous split is inserted:

- (a) A front view of a statue on cement | in a park.
- (b) A front | view of a statue | on cement | in a park.

Our goal is to ensure that every detail unit contains meaningful and novel information, and preserves the grammaticality of the caption. False positives introduce noise that may corrupt the metric signal and compromise training, especially when such errors accumulate.

Given the above reasoning, we modify the segmentation code with several rules to avoid splitting off (1) sentence-initial noun phrases that begin with *The* as they are likely to repeat a previously mentioned entity, (2) prepositional phrases from the noun phrase preceding them as they are likely a modifier to the noun phrase, often referring back to previously mentioned objects (e.g. *The cat is partially* ... in Figure 1), (3) segments which start with a prepositional phrase from the context that follows, unless the segment contains a verb, as they are likely a location modifier to the following noun phrase (e.g. *To the left of the car there is a box*).¹

Negative Triplets For every positive triplet $\{i, c_j, c_j + d_+\}$, we create a negative counterpart, $\{i, c_j, c_j + d_-\}$, by randomly sampling a detail unit from another image-caption pair in the batch. This technique results in what could be called *easy* negatives, i.e. random negatives which can be easily identified as irrelevant to the reference image. Future work could explore the use of hard negatives, but in this work we find that with the right

¹Sometimes, this rule would result in a false negative.

Model	Positive	Negative	Average
CLIP	62.61	68.28	65.44
LongCLIP	60.12	69.93	65.02
SigLIP	58.56	76.28	67.42
NegCLIP	54.96	78.84	66.90
DCI	55.68	63.63	59.66
DAC	46.84	66.88	56.86
LaCLIP	60.98	68.82	64.90
TripletCLIP	53.34	70.20	61.77
LongCLIP*	58.64	77.03	67.83
SPEC	95.37	90.37	92.87

Table 1: Specificity performance of various vision-language models on the sDCI dataset. Positive and Negative correspond to SR_{pos} and SR_{neg} as defined in Section 3.2. *LongCLIP* refers to the ViT-B/16 model as reported in the original paper, while *LongCLIP** is a model we trained from ViT-B/32.

loss weight balancing (see Section 5.4), even this weaker signal can be leveraged effectively.

4.3 Experimental Setup

We train a base LongCLIP-B/32 model with a context window of 248, using standard contrastive training for six epochs. The best checkpoint is then fine-tuned with our specificity objective (see Eq. 4), for another three epochs.

We use the Adam optimizer with a learning rate of 1×10^{-5} , weight decay of 1×10^{-2} , a batch size of 100 per GPU, and gradient accumulation over 4 steps (yielding an effective batch size of 400 per GPU). We set the loss weights to $\alpha = 1$, $\beta = 8$, and $\gamma = 0.8$ based on extensive hyperparameter tuning. All experiments are conducted on four NVIDIA A40 GPUs. Training the model requires approximately one hour per epoch (4 GPU hours).

4.4 Results

Intrinsic Evaluation To evaluate whether our training objective effectively enhances specificity, we compare the specificity rate of SPEC (the base model behind the SPECS metric) against various vision-language models (see Section 3.2) on the sDCI benchmark and report results in Table 1. SPEC achieves the best performance across all VLM models, with $SR_{pos} = 95.37$ and $SR_{neg} = 90.37$, resulting in an average specificity score of 92.87. Compared to the LongCLIP* baseline of 67.83, our model yields a substantial improvement of +25.04 points. The largest gain appears in SR_{pos} ,

where SPEC outperforms LongCLIP* by +36.73, highlighting its superior ability to recognize and prefer more detailed captions, and the effectiveness of the custom training objective.

Interestingly, models with strong general-purpose performance do not necessarily achieve high specificity scores. For example, SigLIP, despite strong results on standard vision-language benchmarks, performs worse than CLIP-based variants in both SR_{pos} and average specificity. This suggests that model architecture alone is not sufficient to capture fine-grained image-text alignment. Models designed to improve compositionality show mixed results: NegCLIP slightly improves over CLIP, while DCI and DAC perform worse, and LaCLIP shows no improvement over the baseline.

Having established that the model we trained shows strong specificity in intrinsic evaluations, we proceed to use this model as a scoring function.

Extrinsic Evaluation To evaluate how well automatic caption metrics align with human preferences, we adopt the evaluation protocol from DE-CAPBENCH (Ye et al., 2025). This human correlation benchmark consists of 100 images, sampled from the ImageInWords (IIW) dataset (Garg et al., 2024). Human-annotated ratings are available for five captions per image, generated by different vision-language models. This setup enables a standardized comparison between automatic metrics and human judgments. We evaluate SPECS in the context of a wide range of metrics from different categories: rule-based, representational similarity-based and LLM-based ones. Table 2 summarizes the results in terms of four standard correlation metrics: Pearson correlation coefficient (PCC), coefficient of determination (R^2) , Kendall's τ (Kd τ), and Sample-wise τ (Sp τ).

Among RS metrics, SPECS achieves the highest human correlation, improving PCC over the CLIPScore baseline from 0.2183 to 0.5228 and Kendall's τ from 0.1724 to 0.4078 . SPECS outperforms most LLM-based metrics, including the strongest open-source metric, FLEUR, in terms of PCC and ranking consistency (Sp τ).

SPECS requires only 2.81×10^{-2} TFLOPs per forward pass, making it far more efficient than LLM-based metrics like FLEUR (7.74) and CLAIR (3.97). With only 0.15 billion parameters, SPECS remains lightweight and scalable, offering a practical and human-aligned solution for evaluating dense, detail-rich captions.

Metric	PCC <i>ρ</i> ↑	1 - R ² ↓	Kd $\tau \uparrow$	$\mathbf{Sp} \ \tau \uparrow$	Base Model	Reference Free	TFLOPs			
Rule-Based Evaluation										
BLEU-4	0.3439	62.78	0.2693	0.2931	_	x	_			
ROUGE	0.2509	156.05	0.1886	0.1893	-	×	_			
METEOR	0.3593	111.95	0.2417	0.2536	_	×	-			
CIDEr	0.0522	3.30E+07	0.0635	0.0601	-	×	-			
	Representational Similarity Evaluation									
SPICE	0.2218	156.11	0.1731	0.1907	-	/	_			
CLIPScore	0.2183	26.04	0.1724	0.1480	CLIP	✓	1.48×10^{-2}			
PACScore	0.1525	20.93	0.1117	0.1260	CLIP	✓	1.48×10^{-2}			
LaCLIP	0.1177	71.94	0.0911	0.1192	CLIP	✓	1.48×10^{-2}			
TripletCLIP	0.1697	34.70	0.0852	0.1038	CLIP	✓	1.48×10^{-2}			
NegCLIP	0.0872	131.57	0.0623	0.0256	CLIP	✓	1.48×10^{-2}			
LongCLIP	0.2320	18.58	0.1769	0.2603	LongCLIP	✓	2.81×10^{-2}			
LongCLIP*	0.1723	33.67	0.1484	0.1662	LongCLIP	✓	2.81×10^{-2}			
SPECS (Ours)	0.5228	3.65	0.4078	0.5400	LongCLIP	1	2.81×10^{-2}			
			LLM-Ba	ased Eval	uation					
FaithScore	0.1937	3.22	0.1626	0.1115	LLaMA	/	3.97			
CLAIR	0.3815	1.98	0.3847	0.4552	Claude	✓	_			
GPT4-Eval	0.3976	2.95	0.3447	0.3866	GPT-4	✓	_			
RLAIF-V	0.3547	5.32	0.2774	0.2544	LLaVA	✓	3.97			
CAPTURE	0.3521	7.62	0.2801	0.3449	InternVL	×	4.54			
FLEUR	0.4230	3.01	0.4246	0.5325	LLaVA	✓	7.37			
DCSCORE	0.6605	1.54	0.5328	0.6166	GPT-40	×	-			

Table 2: Correlation of image captioning evaluation metrics and human judgments: Pearson's ρ (PCC ρ), $1-R^2$, Kendall's τ (Kd τ), and Spearman's τ (Sp τ).For a fair comparison of computational cost, LLM-Based evaluations were computed using an input sequence length of 300 tokens, matching the setting used for Model-Based metrics. Correlation scores lower than those for SPECS are displayed in gray. All p-values are less than 0.001.

5 Further Analysis

5.1 Hubness in the Embedding Space

Although our specificity-enhanced model excels at fine-grained alignment, we observe a decline in performance on standard vision-language tasks such as retrieval and classification. We evaluate generalization across a diverse set of benchmarks, including Urban-1k (Zhang et al., 2024) and COCO (Lin et al., 2014) for text-image retrieval, and ImageNet (Russakovsky et al., 2015), CIFAR-10, and CIFAR-100 (Krizhevsky et al., 2009) for image classification. These benchmarks span both multimodal and unimodal settings, providing a comprehensive view of how specificity-oriented training impacts general-purpose representations (Table 3).

Specifically, our training objective modifies the geometry of the embedding space by introducing additional constraints beyond contrastive similarity, particularly encouraging alignment with incrementally detailed captions. While this improves the model's specificity, it disrupts the isotropy of the

representation space and results in the emergence of hubness: caption embeddings that are overly similar to many images, ultimately degrading retrieval performance.

Overall, while our training strategy enhances specificity evaluation, it can distort the geometry of the embedding space, negatively affecting performance on other downstream tasks. This does not devalue the SPECS metric, but sheds some light into the mechanism it adopts to provide reliable evaluation scores for long image captions.

5.2 Compositionality Analysis

While our main focus is on improving specificity, we also explore whether it leads to improve compositional reasoning. It is reasonable to expect that models capable of handling variations in attribute order or relational structure may also perform better on incrementally positive descriptions. To test this, we evaluate our models on two established benchmarks: ARO (Yuksekgonul et al., 2022), which

Model	Urban-1k		СО	СО	Classification			
	Text-Image	Image-Text	Text-Image	Image-Text	ImageNet	CIFAR-10	CIFAR-100	
CLIP	47.10	61.10	30.45	50.40	68.40	89.75	64.20	
LongCLIP	79.30	79.20	40.40	57.63	66.80	90.69	69.30	
SigLip	62.40	63.10	47.18	65.34	76.08	92.44	72.59	
NegCLIP	52.80	55.60	41.56	56.86	55.84	85.90	60.90	
DCI	43.00	29.70	21.44	20.55	53.34	87.38	57.96	
DAC	23.60	11.40	37.53	33.49	52.36	89.86	64.04	
LongCLIP*	77.00	75.80	35.50	52.44	59.91	90.38	66.36	
SPEC	69.80	0.30	22.72	4.48	11.01	71.26	33.10	

Table 3: Performance of different VL representation models on standard retrieval and classification benchmarks. Results cover long- and short-caption text–image retrieval (Urban-1k, COCO) and image classification (ImageNet, CIFAR-10, CIFAR-100). This table illustrates how specificity-oriented training, while improving fine-grained alignment, can impact general-purpose performance.

Model	Rel.	Attr.	C/O	F/O	SCPP
CLIP	59.84	63.96	47.28	58.54	53.33
LongCLIP	59.70	63.42	56.91	69.03	54.45
SigLip	46.52	56.24	32.95	40.86	20.88
NegCLIP	70.52	81.08	87.04	90.38	63.79
DCI	81.31	73.85	94.53	95.68	51.29
DAC	76.18	67.63	88.58	91.25	43.54
La-CLIP	45.48	58.72	34.97	40.54	54.99
TripletCLIP	54.94	63.07	23.53	27.58	55.71
LongCLIP*	52.96	65.81	63.97	70.20	56.74
SPEC	73.38	69.31	75.23	84.96	35.61

Table 4: Performance of various models on the ARO and SCPP benchmarks. C/O and F/O correspond to compositionality evaluation on COCO-Order and Flickr30k-Order, respectively.

measures understanding of attribute-relation-object structure and word order sensitivity, and Sugar-CREPE++ (SCPP) (Dumpala et al., 2024), which assesses sensitivity to semantic equivalence under lexical variation. Results are shown in Table 4, with full SCPP details in Appendix B.

On ARO, SPEC exhibits considerably higher performance than LongCLIP*, which suggests a direct relationship between specificity and compositionality. This finding does not hold on the SCPP benchmark, however. In fact among all compositionality-enhanced models, only NegCLIP shows a marked improvement on SCPP over the base CLIP model, all others either matching the base performance or showing a considerable degradation (e.g., DAC.)

5.3 Caption Length Sensitivity

We further examined how evaluation performance varies across different caption lengths. To this

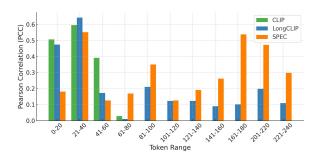


Figure 3: Comparison of Pearson Correlation (PCC) across token ranges for different models on the Flickr30k and ImageInWords datasets.

end, we divided captions into buckets based on token counts and measured human correlation within each range. As shown in Figure 3, CLIPScore performs reasonably well on short captions (fewer than 60 tokens), while SPECS yields consistently stronger correlation for longer captions. This trend suggests that short and long captioning represent two distinct regimes that require different evaluation focuses. Despite attempts to train a unified model that performs well across all lengths, we find a clear trade-off that limits joint performance. We therefore recommend using CLIP for captions under 60 tokens and SPEC for longer ones. The detailed results are provided in Appendix C.

5.4 Hyperparameters

We tuned four key hyperparameters: loss weight (α, β, γ) , learning rate, loss type, and dataset shuffle ratio. Table 5 summarizes the results.

The optimal setting ($\alpha=1, \beta=8, \gamma=0.8$) achieves the highest specificity score of 92.87. Alternative configurations such as 1:9:0.8 and 1:8:0.6 result in noticeably lower performance,

Ablation	Config	Pos.	Neg.	Avg.
Loss Weight 1:8:0.8	1:8:0.6 1:9:0.8 1:8:0.8	85.85 87.39 95.37	82.33 86.59 90.37	84.09 86.99 92.87
Learning Rate 1×10^{-5}	$ \begin{vmatrix} 1 \times 10^{-6} \\ 5 \times 10^{-6} \\ 1 \times 10^{-5} \end{vmatrix} $	77.22 83.57 95.37	68.64 76.54 90.37	72.93 80.55 92.87
Loss Type hinge		80.18 95.37	65.80 90.37	72.99 92.87
Shuffle Rate 90%	50% 100% 90%	83.12 87.27 95.37	87.33 83.11 90.37	85.22 85.19 92.87

Table 5: Hyperparameter tuning.

underscoring the model's sensitivity to the precise relative weighting of different training objectives.

Interestingly, the optimal setting is highly imbalanced, placing much greater emphasis on the positive loss compared to the contrastive and negative loss components. We believe that this imbalance arises from the nature of our specificity-focused training setup: since the contrastive loss is already well optimized from the pretrained CLIP checkpoint, and the negative detail examples are relatively easy, the model benefits more from strong and consistent supervision on the positive detail signal. The positive loss directly encourages the model to increase similarity for incremental, visually grounded additions—precisely the type of finegrained distinction that we aim to capture. Thus, assigning a large weight to this component reinforces the core objective of our method.

We also investigate the role of shuffle ratios when constructing negative captions. Since each negative caption is created by appending a detail unit sourced from other images in the batch, to the current base caption, using unshuffled units may result in semantically coherent and fluent text that unintentionally resembles a valid caption. This risks introducing false negatives that confuse the model during training. To address this, we introduce a shuffle ratio hyperparameter that controls the proportion of negative detail units that are randomly shuffled at the token level before being appended. We find that a shuffle of 90% yields the best performance. A ratio of 90% means most units are shuffled to break semantic coherence, while a small portion (10%) remain in their original order to preserve some challenging cases. This high optimal shuffle rate suggests that introducing controlled

noise into the negatives improves the model's ability to focus on genuine detail alignment without being misled by surface-level fluency.

6 Conclusion

The evaluation of long, detailed captions is a challenge with a pronounced quality-to-cost trade-off. We introduce specificity, a critical dimension for evaluating added detail in image captions. By fine-tuning a CLIP model with a specificity-aware learning objective, we develop SPECS, a new evaluation metric based on representational similarity. Extensive experiments demonstrate that SPECS strongly correlates with human judgments while remaining computationally efficient and scalable.

Limitations

While SPECS offers strong alignment with human judgments and excels at evaluating fine-grained visual grounding, its performance on standard vision-language tasks is limited. As shown in compositionality benchmarks such as ARO and SCPP, improvements in specificity do not directly translate into better reasoning over attribute structures or lexical variations. This indicates that the specificity-focused objective does not generalize well to tasks requiring structural or semantic invariance.

In addition, our hubness analysis reveals distortions in the embedding space caused by specificity-aware training. By encouraging sensitivity to visual details, the model tends to over-align with frequent or stylistically similar captions, leading to degraded performance in retrieval and classification tasks. These findings highlight a trade-off between detail sensitivity and general purpose utility.

Addressing this trade-off remains an open challenge. Future work may consider architectural modifications or auxiliary learning objectives that preserve fine-grained grounding while improving transferability to downstream tasks.

References

Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2016. 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.

Satanjeev Banerjee and Alon Lavie. 2005. 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.
- David Chan, Suzanne Petryk, Joseph E Gonzalez, Trevor Darrell, and John Canny. 2023. Clair: Evaluating image captions with large language models. *arXiv preprint arXiv:2310.12971*.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2024. Sharegpt4v: Improving large multi-modal models with better captions. In *European Conference on Computer Vision*, pages 370–387. Springer.
- Jaemin Cho, Seunghyun Yoon, Ajinkya Kale, Franck Dernoncourt, Trung Bui, and Mohit Bansal. 2022. Fine-grained image captioning with CLIP reward. In Findings of the Association for Computational Linguistics: NAACL 2022, pages 517–527, Seattle, United States. Association for Computational Linguistics.
- Hongyuan Dong, Jiawen Li, Bohong Wu, Jiacong Wang, Yuan Zhang, and Haoyuan Guo. 2024. Benchmarking and improving detail image caption. *arXiv* preprint arXiv:2405.19092.
- Sivan Doveh, Assaf Arbelle, Sivan Harary, Roei Herzig, Donghyun Kim, Paola Cascante-Bonilla, Amit Alfassy, Rameswar Panda, Raja Giryes, Rogerio Feris, and 1 others. 2023. Dense and aligned captions (dac) promote compositional reasoning in vl models. *Advances in Neural Information Processing Systems*, 36:76137–76150.
- Sri Harsha Dumpala, Aman Jaiswal, Chandramouli Shama Sastry, Evangelos Milios, Sageev Oore, and Hassan Sajjad. 2024. Sugarcrepe++ dataset: Visionlanguage model sensitivity to semantic and lexical alterations. *Advances in Neural Information Processing Systems*, 37:17972–18018.
- Lijie Fan, Dilip Krishnan, Phillip Isola, Dina Katabi, and Yonglong Tian. 2023. Improving clip training with language rewrites. *Advances in Neural Information Processing Systems*, 36:35544–35575.
- Roopal Garg, Andrea Burns, Burcu Karagol Ayan, Yonatan Bitton, Ceslee Montgomery, Yasumasa Onoe, Andrew Bunner, Ranjay Krishna, Jason Baldridge, and Radu Soricut. 2024. Imageinwords: Unlocking hyper-detailed image descriptions. *arXiv* preprint arXiv:2405.02793.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. Clipscore: A reference-free evaluation metric for image captioning. *arXiv* preprint arXiv:2104.08718.
- Liqiang Jing, Ruosen Li, Yunmo Chen, Mengzhao Jia, and Xinya Du. 2023. Faithscore: evaluating hallucinations in large vision-language models. corr, abs/2311.01477, 2023. doi: 10.48550. arXiv preprint ARXIV.2311.01477.

- Justin Johnson, Andrej Karpathy, and Li Fei-Fei. 2016. Densecap: Fully convolutional localization networks for dense captioning. In *Proceedings of the IEEE* conference on computer vision and pattern recognition, pages 4565–4574.
- Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, and 1 others. 2023. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. 2009. Cifar-10 and cifar-100 datasets. *URI:* https://www.cs. toronto. edu/kriz/cifar. html, 6(1):1.
- Yebin Lee, Imseong Park, and Myungjoo Kang. 2024. Fleur: An explainable reference-free evaluation metric for image captioning using a large multimodal model. *arXiv preprint arXiv:2406.06004*.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.
- Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. 2021. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In Computer vision—ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13, pages 740–755. Springer.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916.
- Shi Liu, Kecheng Zheng, and Wei Chen. 2024. Paying more attention to image: A training-free method for alleviating hallucination in lvlms. *Preprint*, arXiv:2407.21771.
- Mayu Otani, Riku Togashi, Yu Sawai, Ryosuke Ishigami, Yuta Nakashima, Esa Rahtu, Janne

- Heikkilä, and Shin'ichi Satoh. 2023. Toward verifiable and reproducible human evaluation for text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14277–14286.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. 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.
- Maitreya Patel, Naga Sai Abhiram Kusumba, Sheng Cheng, Changhoon Kim, Tejas Gokhale, Chitta Baral, and 1 others. 2024. Tripletclip: Improving compositional reasoning of clip via synthetic vision-language negatives. *Advances in neural information processing systems*, 37:32731–32760.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, and 1 others. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252.
- Sara Sarto, Manuele Barraco, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. 2023. Positive-augmented contrastive learning for image and video captioning evaluation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6914–6924.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, and 1 others. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294.
- Jack Urbanek, Florian Bordes, Pietro Astolfi, Mary Williamson, Vasu Sharma, and Adriana Romero-Soriano. 2024. A picture is worth more than 77 text tokens: Evaluating clip-style models on dense captions. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 26700–26709.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575.
- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image

- caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3156–3164.
- Yixuan Wang, Qingyan Chen, and Duygu Ataman. 2023. Delving into evaluation metrics for generation: A thorough assessment of how metrics generalize to rephrasing across languages. In *Proceedings of the 4th Workshop on Evaluation and Comparison of NLP Systems*, pages 23–31.
- Chunyu Xie, Bin Wang, Fanjing Kong, Jincheng Li, Dawei Liang, Gengshen Zhang, Dawei Leng, and Yuhui Yin. 2025. Fg-clip: Fine-grained visual and textual alignment. *arXiv preprint arXiv:2505.05071*.
- Zhenlin Xu, Yi Zhu, Siqi Deng, Abhay Mittal, Yanbei Chen, Manchen Wang, Paolo Favaro, Joseph Tighe, and Davide Modolo. 2024. Benchmarking zero-shot recognition with vision-language models: Challenges on granularity and specificity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1827–1836.
- Qinghao Ye, Xianhan Zeng, Fu Li, Chunyuan Li, and Haoqi Fan. 2025. Painting with words: Elevating detailed image captioning with benchmark and alignment learning. *arXiv* preprint arXiv:2503.07906.
- Tianyu Yu, Haoye Zhang, Yuan Yao, Yunkai Dang, Da Chen, Xiaoman Lu, Ganqu Cui, Taiwen He, Zhiyuan Liu, Tat-Seng Chua, and 1 others. 2024. Rlaif-v: Aligning mllms through open-source ai feedback for super gpt-4v trustworthiness. *arXiv preprint arXiv:2405.17220*.
- Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. 2022. When and why vision-language models behave like bags-ofwords, and what to do about it? *arXiv preprint arXiv:2210.01936*.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 11975–11986.
- Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. 2024. Long-clip: Unlocking the long-text capability of clip. In *European Conference on Computer Vision*, pages 310–325. Springer.

A Segmentation Grammar

The following context-free grammar was used to define syntactic structures relevant to caption segmentation:

B SCPP++ Result

Table 7 presents the full results on the SCPP++ benchmark, broken down across five compositional variation types: Swap Object, Swap Attribute, Replace Relation, Replace Object, and Replace Attribute. Each variation is evaluated under two settings: ITT (Image-to-Text retrieval) and TOT (Text-Only Transfer), reflecting different forms of generalization stress.

Overall, we observe that models like NegCLIP and TripletCLIP maintain relatively strong performance across both ITT and TOT settings, while our SPEC model, although competitive in overall specificity evaluation, exhibits lower compositional generalization performance. This is consistent with earlier analysis in Section 5, and supports the claim that specificity-oriented fine-tuning does not necessarily improve compositional reasoning.

C Short Caption Human Correlation

To explore whether a unified evaluation model could perform well across both short and long captions, we evaluated SPECS and CLIP on datasets dominated by shorter captions on short caption datasets. As shown in Table 6, CLIP consistently achieves higher correlation with human judgments across all datasets and metrics. This suggests that while SPECS is optimized for longer, detail-rich captions, it underperforms in short-caption settings. Our results indicate a clear trade-off, and confirm that a single model cannot simultaneously achieve optimal performance across all caption lengths.

$ \begin{array}{c c c c c c c c c c c c c c c c c c c $											
CLIP SPECS 0.6384 2.45 0.4987 0.6226 0.3838 6.78 0.2952 0.4092 COCO Throughness CLIP 0.5785 2.98 0.4458 0.5790 0.5925 0.3645 7.53 0.2784 0.3989 Flickr8k Correctness CLIP 0.5328 3.52 0.4102 0.5422 0.1228 66.7 0.1139 0.1451 Flickr8k Throughness CLIP 0.5012 4.00 0.3790 0.5421 SPECS 0.1995 25.12 0.1550 0.2193 Flickr30k Correctness CLIP 0.6071 2.71 0.4553 0.6219 SPECS 0.2299 18.9 0.1709 0.3058 Flickr30k Throughness CLIP 0.5352 3.49 0.4026 0.5805	Metric	PCC $\rho \uparrow$	$1 - R^2 \downarrow$	Kd $\tau \uparrow$	$ $ Sp $\tau \uparrow$						
CPECS 0.3838 6.78 0.2952 0.4092 COCO Throughness CLIP 0.5785 2.98 0.4458 0.5790 SPECS 0.3645 7.53 0.2784 0.3989 Flickr8k Correctness CLIP 0.5328 3.52 0.4102 0.5422 SPECS 0.1228 66.7 0.1139 0.1451 Flickr8k Throughness CLIP 0.5012 4.00 0.3790 0.5421 SPECS 0.1995 25.12 0.1550 0.2193 Flickr30k Correctness CLIP 0.6071 2.71 0.4553 0.6219 SPECS 0.2299 18.9 0.1709 0.3058 Flickr30k Throughness CLIP 0.5352 3.49 0.4026 0.5805	COCO Correctness										
COCO Throughness	CLIP	0.6384	2.45	0.4987	0.6226						
CLIP 0.5785 2.98 0.4458 0.5790 Flickr8k Correctness CLIP 0.5328 3.52 0.4102 0.5422 SPECS 0.1228 66.7 0.1139 0.1451 Flickr8k Throughness CLIP 0.5012 4.00 0.3790 0.5421 SPECS 0.1995 25.12 0.1550 0.2193 Flickr30k Correctness CLIP 0.6071 2.71 0.4553 0.6219 SPECS 0.2299 18.9 0.1709 0.3058 Flickr30k Throughness CLIP 0.5352 3.49 0.4026 0.5805	SPECS	0.3838	6.78	0.2952	0.4092						
SPECS 0.3645 7.53 0.2784 0.3989 Flickr8k Correctness CLIP 0.5328 3.52 0.4102 0.5422 SPECS 0.1228 66.7 0.1139 0.1451 Flickr8k Throughness CLIP 0.5012 4.00 0.3790 0.5421 SPECS 0.1995 25.12 0.1550 0.2193 Flickr30k Correctness CLIP 0.6071 2.71 0.4553 0.6219 SPECS 0.2299 18.9 0.1709 0.3058 Flickr30k Throughness CLIP 0.5352 3.49 0.4026 0.5805		COCO Throughness									
CLIP 0.5328 3.52 0.4102 0.5422	CLIP	0.5785	2.98	0.4458	0.5790						
CLIP 0.5328 3.52 0.4102 0.5422 SPECS 0.1228 66.7 0.1139 0.1451 Flickr8k Throughness CLIP 0.5012 4.00 0.3790 0.5421 SPECS 0.1995 25.12 0.1550 0.2193 Flickr30k Correctness CLIP 0.6071 2.71 0.4553 0.6219 SPECS 0.2299 18.9 0.1709 0.3058 Flickr30k Throughness CLIP 0.5352 3.49 0.4026 0.5805	SPECS	0.3645	7.53	0.2784	0.3989						
SPECS 0.1228 66.7 0.1139 0.1451 Flickr8k Throughness CLIP 0.5012 4.00 0.3790 0.5421 SPECS 0.1995 25.12 0.1550 0.2193 Flickr30k Correctness CLIP 0.6071 2.71 0.4553 0.6219 SPECS 0.2299 18.9 0.1709 0.3058 Flickr30k Throughness CLIP 0.5352 3.49 0.4026 0.5805		Flickr8k Correctness									
Flickr8k Throughness	CLIP	0.5328	3.52	0.4102	0.5422						
CLIP 0.5012 4.00 0.3790 0.5421 SPECS 0.1995 25.12 0.1550 0.2193 Flickr30k Correctness CLIP 0.6071 2.71 0.4553 0.6219 SPECS 0.2299 18.9 0.1709 0.3058 Flickr30k Throughness CLIP 0.5352 3.49 0.4026 0.5805	SPECS	0.1228	66.7	0.1139	0.1451						
SPECS 0.1995 25.12 0.1550 0.2193 Flickr30k Correctness CLIP 0.6071 2.71 0.4553 0.6219 SPECS 0.2299 18.9 0.1709 0.3058 Flickr30k Throughness CLIP 0.5352 3.49 0.4026 0.5805		Flickr	k Through	hness							
Flickr30k Correctness	CLIP	0.5012	4.00	0.3790	0.5421						
CLIP 0.6071 2.71 0.4553 0.6219 SPECS 0.2299 18.9 0.1709 0.3058 Flickr30k Throughness CLIP 0.5352 3.49 0.4026 0.5805	SPECS	0.1995	25.12	0.1550	0.2193						
SPECS 0.2299 18.9 0.1709 0.3058 Flickr30k Throughness CLIP 0.5352 3.49 0.4026 0.5805		Flickr3	0k Correc	tness							
Flickr30k Throughness CLIP 0.5352 3.49 0.4026 0.5805	CLIP	0.6071	2.71	0.4553	0.6219						
CLIP 0.5352 3.49 0.4026 0.5805	SPECS	0.2299	18.9	0.1709	0.3058						
		Flickr30k Throughness									
SPECS 0.2230 20.12 0.1641 0.2769	CLIP	0.5352	3.49	0.4026	0.5805						
	SPECS	0.2230	20.12	0.1641	0.2769						

Table 6: Human correlation in various short capiton datasets (COCO, Flickr8k, Flickr30k).

D Model Code Names

We provide the exact model code names used in our experiments to ensure reproducibility:

• CLIP: openai/clip-vit-base-patch32

• LongCLIP: BeichenZhang/LongCLIP-B

• SPECS: Xiaohud/SPECS

• FaithScore: llama2/llama-2-7b-hf

• CLAIR: Claude Instant

• GPT-4 Eval: gpt-4-0613

• RLAIF-V: llava-hf/llava-v1.6-7b-hf

CAPTURE: OpenGVLab/InternVL2_5-8B

• FLEUR: llava-hf/llava-v1.6-13b-hf

• **DCSCORE**: gpt-4o-2024-08-06

Model	Swap Object Swap Attribute		Replace Relation		Replace Object		Replace Attribute		Avg.		
	ITT	TOT	ITT	TOT	ITT	TOT	ITT	TOT	ITT	TOT	
CLIP	45.18	19.74	45.21	33.03	56.26	38.62	86.80	83.72	65.61	59.14	53.31
Long-CLIP	42.85	15.10	49.39	31.98	55.68	40.54	90.19	87.71	71.31	59.77	54.45
Long-CLIP*	46.53	28.97	46.99	42.64	52.20	39.68	88.31	91.82	66.37	63.95	56.74
SigLIP	36.32	5.71	30.63	9.00	27.24	12.66	35.16	12.71	30.71	8.62	20.88
NegCLIP	55.25	34.65	57.99	56.47	52.27	51.57	89.53	94.55	69.41	76.27	63.79
DCI	44.10	31.80	45.60	38.00	43.20	35.70	80.20	81.20	60.90	52.20	51.29
DAC	27.80	11.40	33.50	25.40	48.60	48.60	64.30	75.80	44.00	56.00	43.54
La-CLIP	41.22	21.22	48.95	36.04	51.07	42.03	86.44	88.50	68.78	65.61	54.99
TripletCLIP	38.37	18.78	44.44	38.14	58.68	48.08	85.05	89.04	65.61	70.94	55.71
SPEC	30.61	16.73	28.37	24.02	25.96	24.25	48.36	73.91	38.57	45.30	35.61

Table 7: **Compositional Generalization Evaluation.** ITT and TOT denote image-to-text task and text-only task accuracy, respectively.