

# Vision-Language Models Do Not Understand Negation

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

Many practical vision-language applications require models that understand negation, e.g., when using natural language to retrieve images which contain certain objects but not others. Despite advancements in vision-language models (VLMs) through large-scale training, their ability to comprehend negation remains underexplored. This study addresses the question: how well do current VLMs understand negation? We introduce NegBench, a new benchmark designed to evaluate negation understanding across 18 task variations and 79k examples spanning image, video, and medical datasets. The benchmark consists of two core tasks designed to evaluate negation understanding in diverse multimodal settings: Retrieval with Negation and Multiple Choice Questions with Negated Captions. Our evaluation reveals that modern VLMs struggle significantly with negation, often performing at chance level. To address these shortcomings, we explore a data-centric approach wherein we finetune CLIP models on a large-scale synthetic datasets containing millions of negated captions. We show that this approach can result in a 10% increase in recall on negated queries and a 40% boost in accuracy on multiple-choice questions with negated captions.

## 1. Introduction

Joint embedding-based Vision-Language Models (VLMs), such as CLIP, have revolutionized how we approach multi-modal tasks by learning a shared embedding space where both images and text are mapped together. This shared space enables a variety of applications, including cross-modal retrieval, video retrieval, text-to-image generation, image captioning, and even medical diagnosis (CITE[?]). By aligning visual and linguistic representations, these models achieve remarkable performance across domains and are able to model complex interactions between vision and language inputs.

Despite these advances, there is an emerging limitation: these models often fail to handle *negation*, which is essential in many real-world scenarios. Negation allows users

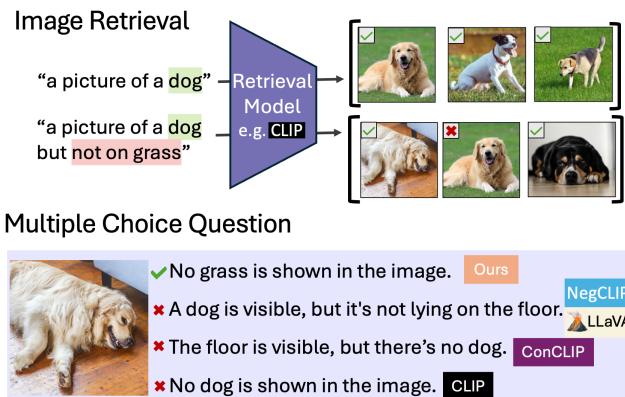


Figure 1. Can Vision-Language Models Understand Negation?

to specify what is absent or false, such as when a radiologist is searching for images showing “bilateral consolidation with no evidence of pneumonia” or a safety inspector querying “construction sites with no barriers”(CITE[?]). Current benchmarks like CREPE and CC-Neg have introduced some forms of negation, but they rely on rigid, templated examples that do not reflect the complexity of natural language queries (CITE[?]). As a result, these benchmarks fall short in evaluating how well VLMs understand negation in practical applications.

To comprehensively evaluate how well VLMs handle negation, we adopt a Two-Stage Evaluation Pipeline. Our pipeline mirrors real-world information retrieval systems, where a coarse-grained retrieval step is followed by a fine-grained ranking or selection step [20, 24].

The first stage, Retrieval-Neg, tests whether models can handle real-world queries that mix affirmative and negative statements, such as “a beach with no people” or “a building without windows.” This task challenges the model to retrieve images from diverse datasets based on the presence of certain elements and the absence of others, simulating scenarios found in search engines, content moderation, and recommendation systems. By retrieving several potentially relevant matches (e.g., top-5 retrieval), Retrieval-Neg serves as the coarse-grained retrieval step in our pipeline.

The second stage, MCQ-Neg, provides a fine-grained,

structured evaluation that directly assesses specific failures in negation. In this task, the model must choose the correct description of an image from several closely related options, where the incorrect choices are hard negatives, differing only by what is affirmed or negated. For instance, in medical diagnostics, consider distinguishing between "The X-ray shows evidence of pneumonia but no evidence of pleural effusion" and "The X-ray shows evidence of pleural effusion but no evidence of pneumonia." These statements are linguistically similar but convey opposite diagnoses, requiring the model to parse subtle yet critical differences.

Through our evaluation pipeline, we uncover a surprising limitation: joint embedding-based VLMs frequently collapse affirmative and negated statements into similar embeddings, treating "a dog" and "no dog" as nearly indistinguishable. This affirmation bias reveals a significant shortcoming that was not sufficiently addressed in previous benchmarks like CREPE or CC-Neg.

Recognizing this critical gap, we then ask: If current models fail to understand negation, can we improve them? To tackle this, we propose a data-centric solution, introducing two large-scale synthetic datasets—Syn-Neg-Cap and Syn-Neg-MCQ—designed to improve negation comprehension. Fine-tuning CLIP on these datasets leads to substantial improvements, including a 10% increase in recall on negated queries and a 40% boost in accuracy on multiple-choice questions with negated captions.

Through NegBench, a systematic error analysis, and a targeted data-driven solution, our work provides a structured approach to understanding and addressing the challenges of negation in VLMs. We will open-source all models and data, hoping to inspire broader research in negation understanding and its wide-ranging applications.

## 2. Related Work

Our work lies within the field of evaluating and advancing foundational vision-language models (VLMs). Joint-embedding models based on CLIP [26] show impressive generalization across visio-linguistic tasks like cross-modal retrieval, image captioning, and visual question answering [2, 15, 16, 25, 27, 30, 33–35] in diverse visual domains, extending beyond natural images to videos and medical images [3, 11, 18, 19, 23, 43]. We introduce a benchmark and data-centric approach to rigorously evaluate and improve negation understanding in these VLMs.

**Negation Understanding in Language and Vision.** Recent work showed that large language models perform sub-optimally when tasked with negation understanding [8]. We go a step further by showing that vision-language models exhibit a more severe affirmation bias, completely failing to differentiate affirmative from negative captions.

Despite this critical limitation, existing benchmarks provide limited assessments of negation in VLMs. CREPE [21]

and the concurrent work CC-Neg [36] are among the few vision-language benchmarks that include negation, but they focus on compositional understanding and rely on linguistic templates that fail to reflect the varied ways negation appears in real user queries. In contrast, our proposed benchmark, NegBench, leverages an LLM to generate natural-sounding negated captions, spanning a broader range of negation types and contexts across images, videos, and medical datasets. This systematic design enables a thorough evaluation of VLMs' ability to handle negation in multimodal settings, uncovering unique challenges and failure cases that have not been fully addressed in prior work.

**Improving CLIP for Compositionality and Negation.** Recent methods have explored improving the generalization abilities of CLIP-like VLMs for visio-linguistic compositionality and limited aspects of negation understanding. For instance, NegCLIP [42] employs composition-aware mining when finetuning CLIP to enhance compositional reasoning, while ConCLIP [36] modifies the CLIP loss to incorporate synthetic, template-based negation examples. In the medical domain, negation is a common feature in clinical text reports, often indicating the absence of specific pathologies [39]. Specialized models like Biomed-CLIP [43] and CONCH [18] have been pretrained on millions of biomedical image-text pairs to address a variety of medical tasks, leveraging domain-specific knowledge from large-scale multimodal data. NegBench provides a systematic way to evaluate general-purpose and medical VLMs.

**Synthetic Data for Model Training.** It is common to use synthetic data to improve the performance of models in computer vision [1, 5, 13, 41]. Recent studies have shown that it is possible to use synthetic data to learn general vision-language representations, with some models trained entirely on synthetic images and captions achieving results comparable to real data [10, 37, 38]. Our approach is similar in spirit, but it constructs synthetic datasets to teach models a new, complex capability—*negation understanding*.

## 3. The Negation Benchmark (NegBench)

We design NegBench as a two-stage pipeline to evaluate the capacity of joint-based vision-language models to understand negation. This structure mirrors real-world applications where an initial, broad retrieval phase is followed by a precise selection or ranking phase.

In the Retrieval-Neg stage, we assess models on their ability to process real-world queries that blend affirmative and negative statements, simulating search tasks where users specify the presence of some elements and the absence of others. The task challenges models to retrieve a set of relevant images (e.g., top-5) that meet these mixed conditions. The second stage, MCQ-Neg, further evaluates models' understanding of negation by having them select the correct description of an image from options that differ only in the

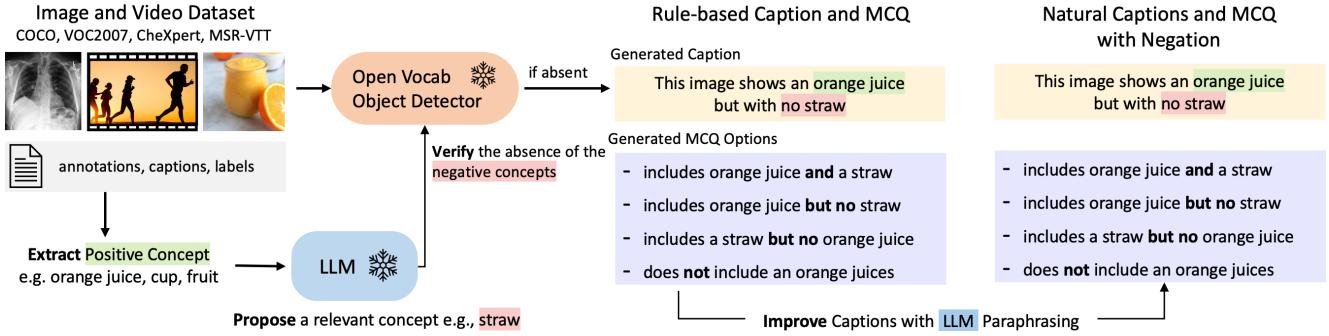


Figure 2. General Pipeline for Constructing NegBench

168

affirmation or negation of specific elements.

169

### 3.1. Transforming Datasets for Negation Evaluation

170

**General Dataset Transformation Overview.** To implement the two-stage evaluation pipeline of NegBench, we adapt several popular vision datasets, covering images (COCO [17], VOC2007), video (MSR-VTT [40]), and specialized medical imaging domains (CheXpert [12]). For each dataset, we identify positive elements  $\{pos\}$ , which represent objects or concepts present in the image, and negative elements  $\{neg\}$ , which are absent from the image but commonly associated with the present objects. When available, we use object-level annotations to identify these elements, as in COCO, VOC2007, and CheXpert; for other datasets, we derive positive and negative elements directly from the captions. This flexible approach allows NegBench to extend any vision dataset, whether it includes object-level annotations or captions, to evaluate negation comprehension across diverse tasks and data modalities.

186

In the Retrieval-Neg task, we modify standard captions by including negations, evaluating how models handle queries that specify both present and absent elements. For example, captions are modified as: “There is no  $x$  in the image. [Original Caption].” or “[Original Caption]. There is no  $x$  in the image.” To introduce linguistic diversity, we use LLaMA 3.1 [6] to paraphrase these captions.

193

For the MCQ-Neg task, we generate multiple-choice questions (MCQs) for each image. The model must identify the correct description based on three linguistic templates: Affirmation, Negation, and Hybrid [14].

197

**1. Affirmation:** “This image includes **A** (and **C**).”

**2. Negation:** “This image does not include **B**.”

**3. Hybrid:** “This image includes **A** but not **B**.”

198

Each MCQ consists of one correct answer and three incorrect answers, which serve as hard negatives, misleading the model if it does not properly understand negation. A correct answer accurately describes the presence of  $\{pos\}$  elements or negates  $\{neg\}$  elements. A False Affirmation (e.g., “This image includes  $x$ ” when  $x \in \{neg\}$ ) or

a False Negation (e.g., “This image does not include  $x$ ” when  $x \in \{pos\}$ ) highlights the model’s failure to comprehend the image. The Hybrid template further evaluates the model’s ability to combine affirmation and negation in the same caption. These MCQs are also paraphrased using LLaMA 3.1 to increase linguistic diversity.

### 3.2. Applicability Across Data Types and Domains

NegBench supports a wide range of data types and domains, enabling comprehensive negation evaluation.

**Video Understanding.** Video retrieval tasks introduce temporal complexity, where negation can involve both objects and actions that vary over time. Using MSR-VTT as an example, we prompt LLaMA 3.1 [6] to extract positive and negative elements from each video’s caption. These elements may represent either objects present in the video or actions taking place. For Retrieval-Neg, we create captions specifying both the presence of some elements and the absence of others (e.g., “A person is cooking but not eating”). In MCQ-Neg, we generate multiple-choice questions where the model must select the description that most accurately represents a video segment, requiring it to reason about negation of objects and actions in dynamic scenes.

**Medical Image Interpretation with CheXpert.** Accurate negation understanding is critical in high-stakes domains like medical imaging. Using the CheXpert dataset [12], we focus on the most frequent condition *Lung Opacity* and design two binary classification tasks:

**Task 1: Affirmation Control Task.** This task evaluates the model’s ability to associate images with specific medical conditions using affirmative statements.

**Question:** Which option describes this image?

A) This image shows Lung Opacity.

B) This image shows Atelectasis.

**Task 2: Negation Understanding Task.** This task tests whether the model can correctly interpret negation, distinguishing the presence or absence of a medical condition.

204

205

206

207

208

209

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

**Question:** Which option best describes the image?

- A) This image shows Lung Opacity.
- B) This image does *not* show Lung Opacity.

238

These extensions highlight the adaptability of NegBench to various data types and domains, from general images and videos to specialized medical imaging. This versatility ensures that NegBench provides rigorous, contextually relevant evaluations of negation understanding in VLMs.

244

### 3.3. Synthetic Datasets for Controlled Evaluation

245  
246

To rigorously test negation understanding, we construct datasets that precisely control object presence and absence.

247

**Motivation and Benefits of Synthetic Data.** Synthetic data offers several advantages over traditional image datasets. First, by creating “hard negatives”—image pairs that differ only by a single object’s presence or absence—we can evaluate the sensitivity of models to negation with minimal confounding variables. Additionally, image datasets like COCO and VOC2007 are limited in the range of visual concepts they cover; COCO has 80 objects while VOC2007 includes only 20. To expand this diversity, we prompt a large language model to propose a broader set of objects, which we use as targets in our synthetic dataset. This approach enables the generation of visually varied scenes that more comprehensively test negation comprehension across a wider array of objects and contexts.

261

**Construction Process for Synthetic Negation Dataset.** We create 10,000 image pairs using Stable Diffusion [29], where each pair includes one image containing a target object and another where it is explicitly absent. To ensure accurate object presence or absence, we use the open-vocabulary object detector OWL-ViT [22].

267

## 4. NegBench Evaluations: Results and Insights

268

In this section, we benchmark the negation abilities of different VLMs using NegBench, comparing models based on their architecture, training data, and training objectives to reveal specific areas where negation understanding remains limited. Specifically, we evaluate five CLIP ViT-B/32 models on Retrieval-Neg and MCQ-Neg tasks. These include OpenAI CLIP [26], CLIP-laion400m [32], and CLIP-datacomp [7], which differ by pretraining dataset, as well as NegCLIP [42], trained to improve compositional language understanding, and ConCLIP [36], trained specifically to improve negation understanding. To handle the video dataset, MSR-VTT, we follow [3] and encode 4 uniformly sampled frames per video, averaging their features to obtain the CLIP video embedding. For medical tasks, we evaluate CONCH [18] and BioMedCLIP [43], two medical

foundation VLMs. We also assess the impact of scaling up CLIP (ViT-B, ViT-L, and ViT-H) to determine if model size improves negation understanding.

283

284

285

**CLIP models struggle with negated queries in retrieval tasks.** We evaluate five CLIP-based models on the original COCO text-to-image retrieval task and its Retrieval-Neg version, where captions include negated statements. Across models, performance drops significantly on the negated task. In COCO retrieval (Figure 3a), CLIP-laion400m experiences a 7.7% drop in recall@5, with CLIP-datacomp and NegCLIP showing drops of 7.6% and 6.8%, respectively. In the more challenging SD-Neg retrieval task (Figure 3b), the performance drops are even more pronounced due to the presence of hard negatives, *i.e.* images that closely resemble positive examples but differ by the exclusion of a single object. Here, NegCLIP, despite its promise for compositional understanding, suffers a 23.0% drop, while ConCLIP, designed specifically for negation understanding, still declines by 18.0%. These results suggest that interpreting negation, particularly in the presence of hard negatives, remains a key challenge for retrieval tasks.

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

**MCQ-Neg reveals severe limitations in CLIP models.** Figure 4a shows that most models perform worse than random guessing (indicated by the red dashed line at 25%) on the MCQ-Neg task, with CLIP-base achieving only 15% on COCO and 8% on VOC2007. These results reveal a fundamental limitation of CLIP’s pretraining objective, which encourages strong associations between visual concepts and specific words, but struggles to interpret negation. Notably, CLIP-laion400m performs better, reaching over 40% accuracy on the SYNTHETIC dataset. This improvement likely stems from the fact that both CLIP-laion400m and Stable Diffusion (used to generate the SYNTHETIC dataset) were trained on the LAION dataset [31]. However, a score of 40% on a 4-way multiple-choice task is still far below an acceptable level, demonstrating that even these models exhibit a serious lack of negation understanding.

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

**Scaling CLIP does not address the negation problem.** As shown in Figure 4b, scaling up the model size from ViT-B/32 (86M parameters) to ViT-L/14 (307M parameters) and ViT-H/14 (632M parameters) does not qualitatively improve negation understanding. While ViT-H/14 performs slightly better on COCO and VOC2007, it underperforms on SD-HNeg and MSR-VTT compared to ViT-B/32. These results suggest that increasing model size alone is not an effective strategy for addressing the fundamental issues with negation understanding.

320

321

322

323

324

325

326

327

328

329

**Critical failures in high-stakes medical tasks.** Figure 4c presents the results for the CheXpert MCQ-Neg task, where

330

331

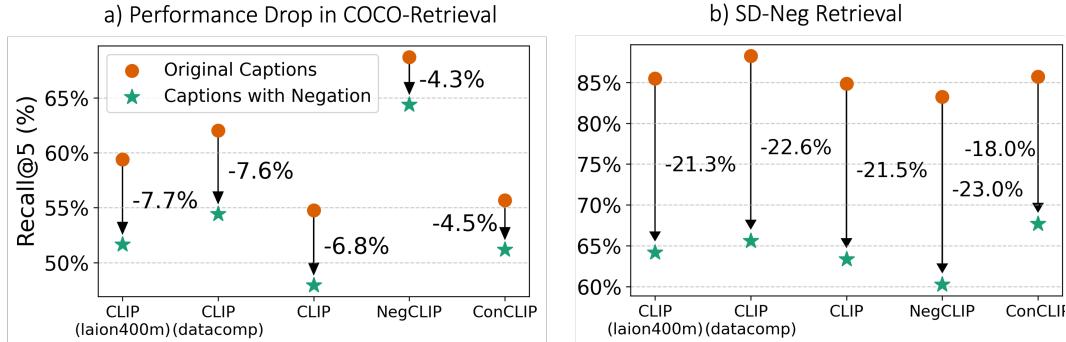


Figure 3. Performance drop in recall@5 on (a) COCO and (b) SD-Neg text-to-image retrieval with negated captions (green stars) compared to original captions (orange circles). All models show substantial drops in performance, with NegCLIP experiencing the largest drop of 23.0% on SD-Neg, which features hard negatives requiring stronger negation reasoning.

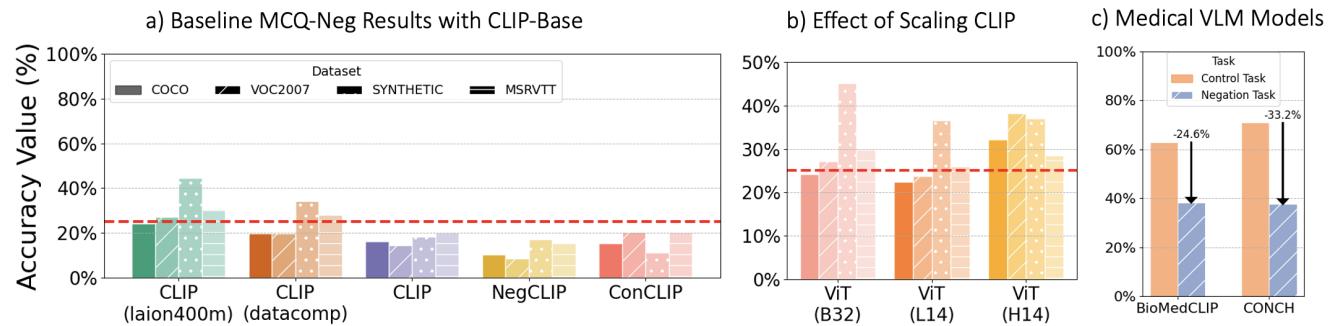


Figure 4. MCQ-Neg performance breakdown for (a) baseline CLIP models, (b) the effect of scaling CLIP, and (c) performance drop of medical VLMs. In (a), CLIP-based models perform worse than random guessing (shown as a red dashed line) on most datasets. In (b), scaling up CLIP models does not significantly improve negation understanding. In (c), medical VLMs models demonstrate a significant drop in performance on negation MCQs.

332 BioMedCLIP and CONCH exhibit substantial performance  
333 drops of 24.6% and 33.2%, respectively, when negation is introduced.  
334 This result is especially concerning in the context of medical diagnostics,  
335 where accurate interpretation of negation (e.g., the presence or absence of a condition such  
336 as Lung Opacity) is essential for correct diagnoses and favorable patient outcomes.  
337  
338

#### 339 4.1. Why Do VLMs Not Understand Negation?

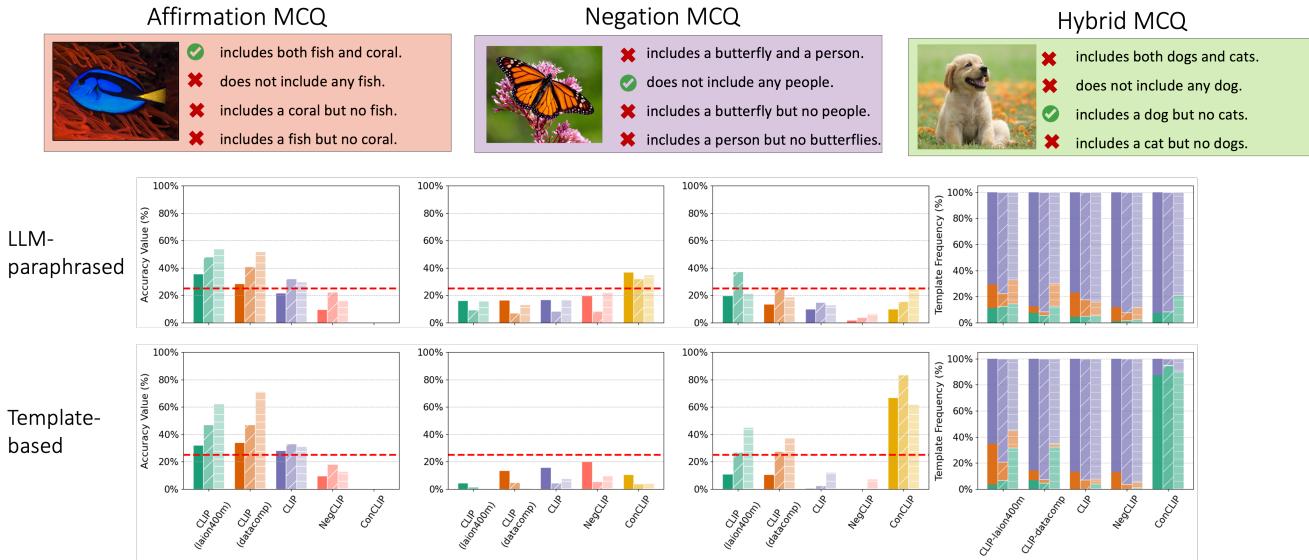
340 The results from NegBench reveal that CLIP VLMs struggle  
341 with different forms of negation understanding, motivating a deeper analysis into the underlying causes of these  
342 failures. In this section, we examine model performance  
343 across different MCQ types and analyze the embedding  
344 spaces of various models to uncover specific shortcut strategies  
345 that limit their negation comprehension.  
346

347 **Model performance varies widely across MCQ types.**  
348 To understand why models perform below random chance,  
349 we categorize the MCQs into three types based on the  
350 correct answer template: Affirmation, Negation, and Hy-

brid. Figure 5 compares model accuracy across these MCQ types, with evaluations conducted in two settings: one using LLaMA 3.1 to paraphrase answer choices into natural-sounding sentences, and another using rigid linguistic templates. All models perform poorly on Negation MCQs, reflecting a general struggle with negation understanding.

351 Most models tend to select Negation sentences regardless  
352 of whether answers are templated or LLM-paraphrased,  
353 as seen in the selection frequencies shown to the right of the  
354 figure. This behavior likely arises from task design, where  
355 67% of MCQs (Negation and Hybrid) lack a correct affirmative  
356 option, leading models to default to “This image  
357 does not include {pos}.” These results suggest that  
358 models trained with CLIP-like objectives often adopt shortcut  
359 strategies that ignore specific words like “no.”

360 The template-based results reveal more biases in model  
361 behavior. For instance, ConCLIP outperforms on Hybrid  
362 MCQs, achieving the highest accuracy, but fails entirely  
363 on Affirmation MCQs, scoring 0% on both image datasets.  
364 This bias is particularly prominent in the rigid template  
365 structure, where ConCLIP is skewed towards constructs like  
366  
367



**Figure 5. Performance breakdown by MCQ type (Affirmation, Negation, Hybrid) across LLM-paraphrased (b) and template-based (c) answer choices.** Models show significant biases towards specific templates. Template selection frequency (right) reveals that models frequently default to Negation answers, especially when a positive object is incorrectly negated.[TODO: Align+Legend]

“This image includes X but not Y.” In fact, as we will show next, ConCLIP maps all templated Hybrid captions to the same location in its embedding space.

**Embedding analysis reveals VLM shortcut strategies.** To investigate potential shortcut strategies, we analyze the embedding spaces of various models using 24 Affirmative (“X”) and 24 Negated (“Not X”) templates to create 48 captions per object. We apply PCA to the resulting embeddings (Figure 6). The templates are included in the appendix.

We observe varying behaviors across models. The overlapping embeddings for affirmative and negated captions in *CLIP* and *NegCLIP* suggest that these models do not distinguish between positive and negative statements, possibly due to a “bag-of-words” shortcut strategy [9, 42] that overlooks negation words. This explains why both models incorrectly select the Negation template, which negates positive objects, in Figure 5. *CoNCLIP* separates positive and negative captions but fails to distinguish between negative captions of different objects, collapsing all negative caption embeddings toward a single point (red circle).

We include the embeddings of a text-only Sentence Transformer [28] as a reference that effectively differentiates affirmative and negated captions along distinct “object type” and “negation” axes, exemplifying ideal separation.

**Hybrid captions reveal more evidence of collapsed embeddings.** Figure 7 extends the previous analysis to hybrid captions that combine affirmations and negations. It provides further evidence that *ConCLIP* employs a shortcut

strategy for embedding linguistic negation, with hybrid and negated captions collapsing towards a single point (green circle), indicating significant compression along the negation axis. While *CLIP* and *NegCLIP* struggle to distinguish affirmative from negative statements, *NegCLIP* shows better separation for hybrid captions, which appear collapsed in the *CLIP* embedding space. This suggests that *NegCLIP*’s poor performance on Hybrid MCQs might be due to a misalignment between the text and image encoders, rather than an inability to understand hybrid sentence structure. In contrast, the *Sentence Transformer* effectively distinguishes between different caption types and provides semantically guided representations. For example, it aligns “flowers but not cats” along the line connecting “cats” and “not flowers.”

## 5. A Data-Centric Approach for Improving Negation Understanding

Our analysis in Section 4.1 suggests that the poor negation understanding of CLIP-based models is caused by overreliance on linguistic shortcut strategies. We hypothesize that these shortcuts arise from data limitations during training. In CLIP, training data lacks examples with explicit negation, leaving it unable to distinguish negated and affirmed concepts. In contrast, ConCLIP’s training data overfits to a single hybrid linguistic template, limiting its ability to generalize across varied negation structures. In this section, we explore data-centric strategies to address these gaps, introducing a dataset that not only includes diverse negation examples but also spans a range of linguistic styles, aiming to

400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413

414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427

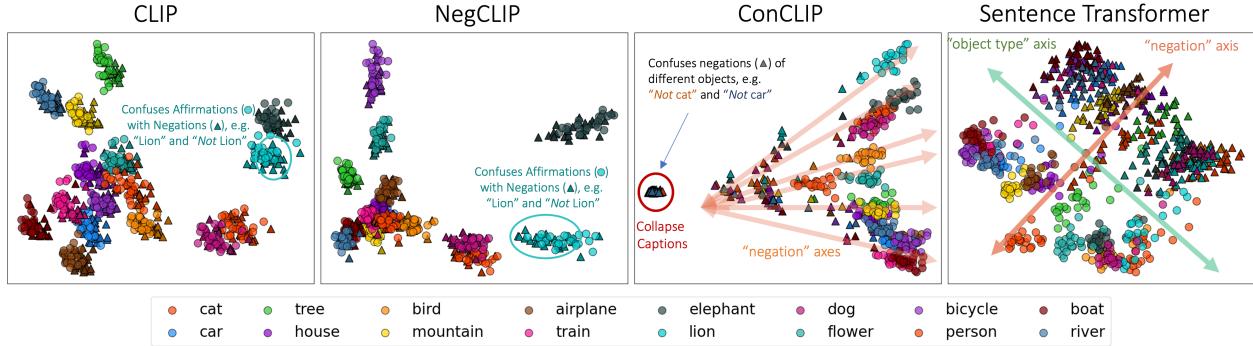


Figure 6. **PCA projections of Affirmative (dots) and Negative (triangles) Caption Embeddings for Multiple Objects.** The Sentence Transformer shows clear separation along 'object type' and 'negation' axes, while CoNCLIP demonstrates compression along negation axes (treating all negated captions as identical). CLIP and NegCLIP lack separation between affirmative and negated captions.

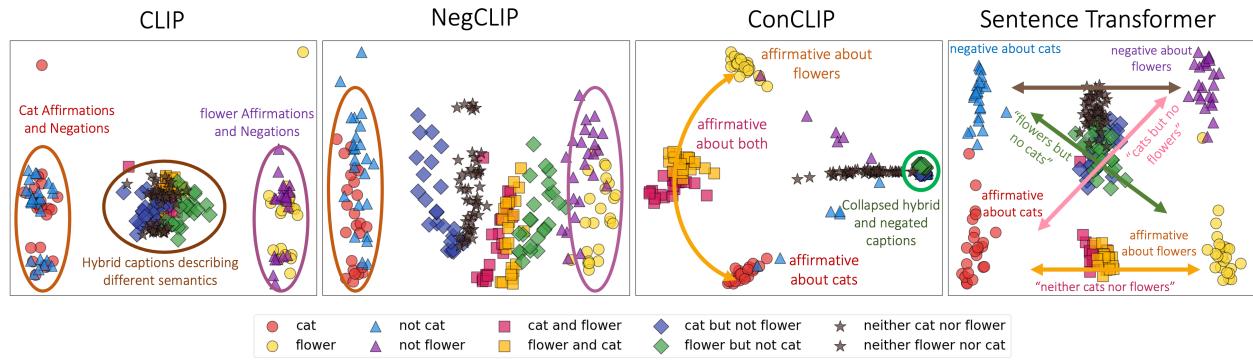


Figure 7. **PCA of Affirmative (dots and squares), Negated (triangles and stars), Hybrid captions (diamonds).** ConCLIP compresses negated and hybrid captions (green circle), while CLIP and NegCLIP fail to separate affirmative and negated captions (red and purple ovals). Sentence Transformer exhibits clear separability and semantic representation, exemplifying the desired behavior.

428 aiming to create embedding spaces that effectively separate  
429 affirmative and negated concepts.

### 430 5.1. Synthetic Negation Dataset Creation

431 We augment the CC12M dataset [4], which contains ap-  
432 proximately 10 million image-text pairs, to generate two  
433 synthetic datasets with negated captions: Syn-Neg-Cap and  
434 Syn-Neg-MCQ. Our goal is to expose models to a wide vari-  
435 ety of negation scenarios and improve their ability to encode  
436 negated statements. The process follows these steps:

437 1. **Object Extraction:** Using LLaMA 3.1 [6], we extract  
438 positive objects (those mentioned in the caption) and  
439 negative objects (contextually relevant but not present)  
440 from each image-caption pair in CC12M.

441 2. **Visual Verification:** An open-vocabulary object detec-  
442 tor [22] verifies the presence of positive objects and en-  
443 sures the absence of the negative objects in the image.  
444 This step is crucial to avoid introducing incorrect nega-  
445 tions that could confuse the model.

446 3. **Caption Generation:** For each image, we generate mul-  
447 tiple new captions that incorporate negated objects into

the original captions. LLaMA 3.1 is used to ensure the generated captions are natural-sounding and reflect realistic negation scenarios found in retrieval queries. This process is flexible and agnostic to the exact number of captions generated per image.

448 We construct two variants of the synthetic dataset:

- 449 • **Syn-Neg-Cap:** For each image in CC12M, we generate  
450 three captions incorporating negated objects, resulting in  
451 approximately 30 million synthetic captions.
- 452 • **Syn-Neg-MCQ:** For this variant, we generate four cap-  
453 tions per image, one of which is correct and the other  
454 three are hard negatives, selected based on object anno-  
455 tations. This format provides stronger training signals  
456 for fine-grained negation understanding. Syn-Neg-MCQ  
457 consists of around 40 million synthetic captions.

458 We will release the extracted object annotations for each  
459 image in CC12M, along with the corresponding URLs, and  
460 all the generated captions in Syn-Neg-Cap and Syn-Neg-  
461 MCQ. This will help the community build on our dataset  
462 and advance negation understanding and multimodal re-  
463 trieval research.

## 469 5.2. Fine-Tuning with Negation-Enriched Data

470 **Standard CLIP Objective on Syn-Neg-Cap.** Let  $\mathcal{B}_{\text{cap}} = \{(I_i, T_i)\}_{i=1}^N$  represent a batch of  $N$  image-caption pairs  
 471 from Syn-Neg-Cap, where each image  $I_i$  is paired with a  
 472 caption  $T_i$  that describes present and absent objects in the  
 473 image. For each batch  $\mathcal{B}_{\text{cap}}$ , we compute a similarity matrix  
 474  $S \in \mathbb{R}^{N \times N}$ , where each element  $S_{j,k}$  represents the cosine  
 475 similarity between the  $j$ -th image and the  $k$ -th caption. The  
 476 CLIP objective applies a symmetric cross-entropy loss over  
 477 this matrix, encouraging high similarity for correct image-  
 478 caption pairs and low similarity for incorrect pairs. This  
 479 loss is denoted as  $\mathcal{L}_{\text{CLIP}}(\mathcal{B}_{\text{cap}})$  and provides the model with  
 480 diverse negation examples in a contrastive learning setup.  
 481

482 **Multiple-Choice Objective on Syn-Neg-MCQ.** Let  
 483  $\mathcal{B}_{\text{mcq}} = \{(I_i, T_{i,1}, T_{i,2}, T_{i,3}, T_{i,4})\}_{i=1}^M$  represent a batch of  
 484  $M$  examples from Syn-Neg-MCQ, where each image  $I_i$  has  
 485 four associated captions  $\{T_{i,j}\}_{j=1}^4$ , with one caption cor-  
 486 rectly describing the image and the others serving as hard  
 487 negatives. To fine-tune on Syn-Neg-MCQ, we compute the  
 488 cosine similarity between each image and its four caption  
 489 options to obtain a set of logits for each image-option pair.  
 490

491 The multiple-choice loss  $\mathcal{L}_{\text{MCQ}}(\mathcal{B}_{\text{mcq}})$  is then computed  
 492 by applying a cross-entropy loss over the logits, with the  
 493 correct answer index as the target. This loss encourages the  
 494 model to assign higher similarity to the correct caption and  
 lower similarity to the hard negative captions:

$$495 \mathcal{L}_{\text{MCQ}}(\mathcal{B}_{\text{mcq}}) = -\frac{1}{M} \sum_{i=1}^M \log \frac{\exp(\text{logits}_{i,\text{correct\_answers}_i})}{\sum_{j=1}^4 \exp(\text{logits}_{i,j})}, \quad (1)$$

496 where  $\text{correct\_answers}_i$  indicates the index of the correct  
 497 caption for the  $i$ -th image.

498 **Combined Training Objective.** The final objective com-  
 499 bines the contrastive loss on Syn-Neg-Cap with the MCQ  
 500 loss on Syn-Neg-MCQ, weighted by  $\alpha$  and  $\beta$  to balance  
 501 their contributions. The total loss for one batch is:

$$502 \mathcal{L}_{\text{Total}} = \alpha \mathcal{L}_{\text{CLIP}}(\mathcal{B}_{\text{cap}}) + \beta \mathcal{L}_{\text{MCQ}}(\mathcal{B}_{\text{mcq}}). \quad (2)$$

503 **Evaluation Protocol.** To assess the impact of our data-  
 504 centric approach, we fine-tune two pretrained models (Ope-  
 505 nAI CLIP and NegCLIP) on Syn-Neg-Cap using the con-  
 506 trastive loss  $\mathcal{L}_{\text{CLIP}}$ . Additionally, we fine-tune both models  
 507 on the combined Syn-Neg-Cap and Syn-Neg-MCQ datasets  
 508 using  $\mathcal{L}_{\text{Total}}$  in Equation (2). For comparison, we fine-tune  
 509 these models on the original CC12M dataset to isolate the  
 510 effect of our negation-enriched datasets. Our goal is to  
 511 demonstrate that CLIP models can significantly improve  
 512 their understanding of negation with the right data.

Model	Fine-tune data	R@5 ( $\uparrow$ )	R-Neg@5 ( $\uparrow$ )	MCQ ( $\uparrow$ )
CLIP	None	54.8	48.0	16.3
	CC12M	58.8	54.5	11.2 ( $\downarrow 5.1$ )
	<b>Syn-Neg-Cap</b>	<b>58.5</b>	<b>57.8</b>	<b>14.7</b> ( $\downarrow 1.6$ )
	<b>Syn-Neg-Full</b>	<b>54.2</b>	<b>51.9</b>	<b>46.9</b> ( $\uparrow 30.6$ )
NegCLIP	None	68.7	64.4	10.2
	CC12M	70.2	66.0	10.6 ( $\uparrow 0.4$ )
	<b>Syn-Neg-Cap</b>	<b>68.6</b>	<b>67.5</b>	<b>12.5</b> ( $\uparrow 2.3$ )
	<b>Syn-Neg-Full</b>	<b>69.0</b>	<b>67.0</b>	<b>51.0</b> ( $\uparrow 40.8$ )

(a) COCO Evaluation

Model	Fine-tune data	R@5 ( $\uparrow$ )	R-Neg@5 ( $\uparrow$ )	MCQ ( $\uparrow$ )
CLIP	None	50.6	45.8	20.1
	CC12M	53.7	49.9	16.9 ( $\downarrow 3.2$ )
	<b>Syn-Neg-Cap</b>	<b>54.1</b>	<b>53.5</b>	<b>20.1</b> (0.0)
	<b>Syn-Neg-Full</b>	<b>46.9</b>	<b>43.9</b>	<b>35.6</b> ( $\uparrow 15.5$ )
NegCLIP	None	53.7	51.0	15.3
	CC12M	56.4	52.6	16.8 ( $\uparrow 1.5$ )
	<b>Syn-Neg-Cap</b>	<b>56.5</b>	<b>54.6</b>	<b>18.9</b> ( $\uparrow 3.6$ )
	<b>Syn-Neg-Full</b>	<b>54</b>	<b>51.5</b>	<b>36.6</b> ( $\uparrow 21.3$ )

(b) MSR-VTT Evaluation

Table 1. **Comparison of fine-tuning datasets** on performance metrics across COCO and MSR-VTT, fine-tuned on respective datasets and evaluated on retrieval and multiple-choice questions. Differences in MCQ accuracy from the baseline are shown, with increases of  $\uparrow 1$  or more highlighted.

We evaluate the models on two tasks: (i) text-to-image and text-to-video retrieval on COCO and MSR-VTT, both with and without negated queries, and (ii) image-to-text and video-to-text MCQ tasks, where models select the correct caption from four options. The results are shown in Table 1.

**Results.** Fine-tuning CLIP on Syn-Neg-Cap leads to significant improvements in handling negated queries in retrieval tasks. On COCO, the R-Neg@5 score increases by 10.6%, while the gap between R@5 and R-Neg@5 ( $\Delta R$ ) narrows from 6.8% to 0.5%, indicating that the model now performs nearly as well on negated queries as on standard ones. A similar pattern is observed on MSR-VTT.

However, fine-tuning on Syn-Neg-Cap alone does not improve performance on the MCQ task, suggesting that the contrastive objective is insufficient for learning fine-grained negation comprehension. To address this, we fine-tune CLIP and NegCLIP on the combined Syn-Neg-Cap and Syn-Neg-MCQ datasets, yielding substantial improvements on MCQ tasks. On COCO-MCQ, for instance, CLIP’s accuracy rises from 16.3% to 56.9%, a 40.6% increase.

Overall, fine-tuning on our synthetic datasets significantly enhances the models’ negation understanding without sacrificing performance on standard retrieval tasks. This demonstrates that CLIP-based models can be made more robust to negated queries with the right data, reducing their reliance on shortcut strategies.

**Ablation: Effect of varying  $\alpha$ .** The table below shows the impact of varying the weight factor  $\alpha$  in the combined

loss  $\mathcal{L}_{\text{Total}} = \alpha \mathcal{L}_{\text{CLIP}} + (1 - \alpha) \mathcal{L}_{\text{MCQ}}$  when fine-tuning CLIP on Syn-Neg-Full. As  $\beta$  increases, more weight is placed on the original CLIP contrastive objective, while a lower  $\beta$  emphasizes the MCQ loss.

$\alpha$	0	0.5	0.9	0.99	1
COCO Recall@5 (%)	33.9	37.3	47.6	54.2	58.5
COCO MCQ Acc (%)	61.0	54.7	50.5	46.9	14.7

Properly tuning  $\alpha$  is important to strike a balance, adapting to fine-grained MCQs while maintaining standard retrieval performance.

## References

- [1] Hassan Abu Alhaija, Siva Karthik Mustikovela, Lars Mescheder, Andreas Geiger, and Carsten Rother. Augmented reality meets computer vision: Efficient data generation for urban driving scenes. *International Journal of Computer Vision*, 126:961–972, 2018. [2](#)
- [2] Alberto Baldrati, Marco Bertini, Tiberio Uricchio, and Alberto Del Bimbo. Effective conditioned and composed image retrieval combining clip-based features. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 21466–21474, 2022. [2](#)
- [3] Santiago Castro and Fabian Caba. FitCLIP: Refining large-scale pretrained image-text models for zero-shot video understanding tasks. In *33rd British Machine Vision Conference 2022, BMVC 2022, London, UK, November 21-24, 2022*. BMVA Press, 2022. [2, 4](#)
- [4] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021. [7](#)
- [5] Yuhua Chen, Wen Li, Xiaoran Chen, and Luc Van Gool. Learning semantic segmentation from synthetic data: A geometrically guided input-output adaptation approach. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1841–1850, 2019. [2](#)
- [6] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. [3, 7](#)
- [7] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, Eyal Orgad, Rahim Entezari, Giannis Daras, Sarah M Pratt, Vivek Ramanujan, Yonatan Bitton, Kalyani Marathe, Stephen Mussmann, Richard Vencu, Mehdi Cherti, Ranjay Krishna, Pang Wei Koh, Olga Saukh, Alexander Ratner, Shuran Song, Hannaneh Hajishirzi, Ali Farhadi, Romain Beaumont, Sewoong Oh, Alex Dimakis, Jenia Jitsev, Yair Carmon, Vaishaal Shankar, and Ludwig Schmidt. Datacomp: In search of the next generation of multimodal datasets. In *NeurIPS Datasets and Benchmarks Track*, 2023. [4](#)
- [8] Iker García-Ferrero, Begoña Altuna, Javier Alvez, Itziar Gonzalez-Dios, and German Rigau. This is not a dataset: A large negation benchmark to challenge large language models. In *EMNLP*. Association for Computational Linguistics, 2023. [2](#)
- [9] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11), 2020. [6](#)
- [10] Hasan Abed Al Kader Hammoud, Hani Itani, Fabio Pizzati, Philip Torr, Adel Bibi, and Bernard Ghanem. Synthclip: Are we ready for a fully synthetic clip training? *arXiv preprint arXiv:2402.01832*, 2024. [2](#)
- [11] Wisdom Ikezogwo, Saygin Seyfioglu, Fatemeh Ghezloo, Dylan Geva, Fatwir Sheikh Mohammed, Pavan Kumar Anand, Ranjay Krishna, and Linda Shapiro. Quilt-1m: One million image-text pairs for histopathology. *Advances in neural information processing systems*, 36, 2024. [2](#)
- [12] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoor, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence*, pages 590–597, 2019. [3](#)
- [13] Ali Jahanian, Xavier Puig, Yonglong Tian, and Phillip Isola. Generative models as a data source for multiview representation learning. In *International Conference on Learning Representations*, 2022. [2](#)
- [14] Miren Itziar Laka Mugarza. *Negation in syntax—on the nature of functional categories and projections*. PhD thesis, Massachusetts Institute of Technology, 1990. [3](#)
- [15] Yehao Li, Yingwei Pan, Ting Yao, and Tao Mei. Comprehending and ordering semantics for image captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 17990–17999, 2022. [2](#)
- [16] Zhengxin Li, Wenzhe Zhao, Xuanyi Du, Guangyao Zhou, and Songlin Zhang. Cross-modal retrieval and semantic refinement for remote sensing image captioning. *Remote Sensing*, 16(1):196, 2024. [2](#)
- [17] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. [3](#)
- [18] Ming Y Lu, Bowen Chen, Drew FK Williamson, Richard J Chen, Ivy Liang, Tong Ding, Guillaume Jaume, Igor Odintsov, Long Phi Le, Georg Gerber, et al. A visual-language foundation model for computational pathology. *Nature Medicine*, 30:863–874, 2024. [2, 4](#)
- [19] Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. CLIP4Clip: An empirical study of clip for end to end video clip retrieval. *arXiv preprint arXiv:2104.08860*, 2021. [2](#)
- [20] Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. Fine-tuning llama for multi-stage text retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2421–2425, 2024. [1](#)
- [21] Zixian Ma, Jerry Hong, Mustafa Omer Gul, Mona Gandhi, Irena Gao, and Ranjay Krishna. Crepe: Can vision-language

- 651 foundation models reason compositionally? In *CVPR*, 2023. 2
- 652
- 653 [22] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim 708  
654 Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh 709  
655 Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran 710  
656 Shen, Xiao Wang, Xiaohua Zhai, Thomas Kipf, and Neil 711  
657 Houlsby. Simple open-vocabulary object detection with 712  
658 vision transformers. *ECCV*, 2022. 4, 7
- 659 [23] Medhini Narasimhan, Anna Rohrbach, and Trevor Darrell. 713  
660 Clip-it! language-guided video summarization. *Advances 714  
661 in neural information processing systems*, 34:13988–14000, 715  
662 2021. 2
- 663 [24] Rodrigo Nogueira, Wei Yang, Kyunghyun Cho, and Jimmy 716  
664 Lin. Multi-stage document ranking with bert. *arXiv preprint 717  
665 arXiv:1910.14424*, 2019. 1
- 666 [25] Suzanne Petryk, Lisa Dunlap, Keyan Nasseri, Joseph 718  
667 Gonzalez, Trevor Darrell, and Anna Rohrbach. On guiding 719  
668 visual attention with language specification. In *Proceedings 720  
669 of the IEEE/CVF Conference on Computer Vision and Pattern 721  
670 Recognition*, pages 18092–18102, 2022. 2
- 671 [26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya 722  
672 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, 723  
673 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning 724  
674 transferable visual models from natural language super- 725  
675 vision. In *ICML*. PMLR, 2021. 2, 4
- 676 [27] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong 726  
677 Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. 727  
678 Denseclip: Language-guided dense prediction with context- 728  
679 aware prompting. In *Proceedings of the IEEE/CVF conference 729  
680 on computer vision and pattern recognition*, pages 730  
681 18082–18091, 2022. 2
- 682 [28] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence 731  
683 embeddings using siamese bert-networks. In *EMNLP*. Association 732  
684 for Computational Linguistics, 2019. 6
- 685 [29] Robin Rombach, Andreas Blattmann, Dominik Lorenz, 733  
686 Patrick Esser, and Björn Ommer. High-resolution image 734  
687 synthesis with latent diffusion models. In *Conference on 735  
688 Computer Vision and Pattern Recognition (CVPR)*, pages 10684– 736  
689 10695, 2022. 4
- 690 [30] Aneeshan Sain, Ayan Kumar Bhunia, Pinaki Nath Chowdhury, 737  
691 Subhadeep Koley, Tao Xiang, and Yi-Zhe Song. Clip 738  
692 for all things zero-shot sketch-based image retrieval, fine- 739  
693 grained or not. In *Proceedings of the IEEE/CVF Conference 740  
694 on Computer Vision and Pattern Recognition*, pages 2765– 741  
695 2775, 2023. 2
- 696 [31] Christoph Schuhmann, Romain Beaumont, Richard Vencu, 742  
697 Cade Gordon, Ross Wightman, Mehdi Cherti, Theo 743  
698 Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, 744  
699 et al. Laion-5b: An open large-scale dataset for training 745  
700 next generation image-text models. *Advances in Neural 746  
701 Information Processing Systems*, 35:25278–25294, 2022. 4
- 702 [32] Christoph Schuhmann, Romain Beaumont, Richard Vencu, 747  
703 Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo 748  
704 Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, 749  
705 Patrick Schramowski, Srivatsa R Kundurthy, Katherine 750  
706 Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia 751  
707 Jitsev. LAION-5b: An open large-scale dataset for train- 752  
708 ing next generation image-text models. In *Thirty-sixth Con- 709  
709 ference on Neural Information Processing Systems Datasets 710  
710 and Benchmarks Track*, 2022. 4
- 711 [33] Sheng Shen, Liunian Harold Li, Hao Tan, Mohit Bansal, 712  
712 Anna Rohrbach, Kai-Wei Chang, Zhewei Yao, and Kurt 713  
713 Keutzer. How much can clip benefit vision-and-language 714  
714 tasks? In *International Conference on Learning Representa- 715  
715 tions*. 2
- 716 [34] Hengcan Shi, Munawar Hayat, Yicheng Wu, and Jianfei 717  
717 Cai. Proposalsclip: Unsupervised open-category object 718  
718 proposal generation via exploiting clip cues. In *Proceedings 719  
719 of the IEEE/CVF Conference on Computer Vision and Pattern 720  
720 Recognition*, pages 9611–9620, 2022.
- 721 [35] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: 722  
722 What and where pathways for robotic manipulation. In *Con- 723  
723 ference on robot learning*, pages 894–906. PMLR, 2022. 2
- 724 [36] Jaisidh Singh, Ishaaan Shrivastava, Mayank Vatsa, Richa 725  
725 Singh, and Aparna Bharati. Learn “no” to say “yes” bet- 726  
726 ter: Improving vision-language models via negations. *arXiv 727  
727 preprint arXiv:2403.20312*, 2024. 2, 4
- 728 [37] Yonglong Tian, Lijie Fan, Phillip Isola, Huiwen Chang, and 729  
729 Dilip Krishnan. Stablerep: Synthetic images from text-to- 730  
730 image models make strong visual representation learners. In 731  
731 *NeurIPS*, 2023. 2
- 732 [38] Yonglong Tian, Lijie Fan, Kaifeng Chen, Dina Katabi, Dilip 733  
733 Krishnan, and Phillip Isola. Learning vision from mod- 734  
734 els rivals learning vision from data. In *Proceedings of 735  
735 the IEEE/CVF Conference on Computer Vision and Pattern 736  
736 Recognition*, pages 15887–15898, 2024. 2
- 737 [39] Ekin Tiu, Ellie Talius, Pujan Patel, Curtis P Langlotz, Andrew 738  
738 Y Ng, and Pranav Rajpurkar. Expert-level detection 739  
739 of pathologies from unannotated chest x-ray images via self- 740  
740 supervised learning. *Nature Biomedical Engineering*, 6(12): 741  
1399–1406, 2022. 2
- 742 [40] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large 743  
743 video description dataset for bridging video and language. In 744  
744 *Proceedings of the IEEE conference on computer vision and 745  
745 pattern recognition*, pages 5288–5296, 2016. 3
- 746 [41] Jianhao Yuan, Jie Zhang, Shuyang Sun, Philip Torr, and Bo 747  
747 Zhao. Real-fake: Effective training data synthesis through 748  
748 distribution matching. In *The Twelfth International Conference 749  
749 on Learning Representations*, 2024. 2
- 750 [42] Mert Yuksekogul, Federico Bianchi, Pratyusha Kalluri, 751  
751 Dan Jurafsky, and James Zou. When and why vision- 752  
752 language models behave like bags-of-words, and what to do 753  
753 about it? In *ICLR*, 2023. 2, 4, 6
- 754 [43] Sheng Zhang, Yanbo Xu, Naoto Usuyama, Hanwen Xu, 755  
755 Jaspreet Bagga, Robert Tinn, Sam Preston, Rajesh Rao, 756  
756 Mu Wei, Naveen Valluri, et al. Biomedclip: a multimodal 757  
757 biomedical foundation model pretrained from fifteen million 758  
758 scientific image-text pairs. *arXiv preprint arXiv:2303.00915*, 759  
759 2023. 2, 4