TNG-CLIP: Training-Time Negation Data Generation for Negation Awareness of CLIP

Anonymous ACL submission

Abstract

Vision-language models (VLMs), such as CLIP, have demonstrated strong performance across a range of downstream tasks. However, CLIP is still limited in negation understanding: the ability to recognize the absence or exclusion of a concept. Existing methods address the problem by using a large language model (LLM) to generate large-scale data of image captions containing negation for further fine-tuning CLIP. However, these methods are both time- and computeintensive, and their evaluations are typically restricted to image-text matching tasks. To expand the horizon, we (1) introduce a trainingtime negation data generation pipeline such that negation captions are generated during the training stage, which only increases 2.5% extra training time, and (2) we propose the first benchmark, NEG-TTOI, for evaluating text-toimage generation models on prompts containing negation, assessing model's ability to produce semantically accurate images. We show that our proposed method, TNG-CLIP, achieves SOTA performance on diverse negation benchmarks of image-to-text matching, text-to-image retrieval, and image generation.

1 Introduction

007

011

014 015

017

037

041

Vision-language models (VLM), such as CLIP (Radford et al., 2021), provide an efficient approach to tackle vision-language tasks by learning the features of different modalities in a shared embedding space. However, these models fundamentally lack a robust understanding of **negation**—the ability to recognize the absence or exclusion of a concept, *e.g.*, "A dog **not** playing a ball.", "There is **no** tree on the street.". Negation is a fundamental aspect of human reasoning, enabling precise description of constraints and expectations in communication. Without proper negation understanding, VLMs generate and retrieve semantically incorrect content, particularly in complicated scenarios where the



Figure 1: We present *TNG-CLIP*, a negation-aware CLIP that achieves outstanding negation understanding in image-to-text matching, text-to-image retrieval and proposed image generation NEG-TTOI benchmarks.

presence or absence of specific elements critically alters meanings.

To tackle this problem, current methods (Alhamoud et al., 2025a; Singh et al., 2024; Park et al., 2025; Yuksekgonul et al., 2023) focus on generating well-designed image-text datasets, such that there are negation captions associated with each image sample, and then fine-tune the underlying VLM. However, such approaches face three challenges: (1) the negation of each caption is designed, generated, and verified via LLMs. Considering the fact that the existing vision-language datasets (Chen et al., 2015; Changpinyo et al., 2021) contain millions of samples, generating the negation dataset is extremely time- and compute-consuming. (2) Unlike standard semantic descriptions, which are typically grounded in observable features, the negation process introduces arbitrariness by specifying the absence of concepts that are not depicted. For example, given an image of "a dog playing a ball", one could construct multiple valid negation captions such as "a dog playing a ball while no man is present" or "a dog playing a ball but not on the beach". By generating fixed negation captions, previous methods may constrain the diversity of negation scenarios, thus harming the generalization

043

044

068of the fine-tuned VLM on negation understanding069tasks. (3) Previous methods are mainly evaluated070on image-to-text matching and text-to-image re-071trieval tasks. Considering the versatility of CLIP,072however, evaluation should not be constrained to073matching-based tasks and must include more di-074verse downstream tasks such as generation-based075tasks, where the text encoder can be used as part of076a generative model (Rombach et al., 2022).

079

097

100

101

102

103

104

106

108

109

110

111

112

113

114

115

116

117

118

We propose a new data generation and training pipeline which generates negation captions during training without the need for a pre-defined negated image-text pair dataset. In each training batch, we identify the most similar image-text pair for every image-text example by computing the cosine similarity between their embedded image features. For each caption, we generate negated variants using a template-based approach, by interacting with another caption in the same batch. Because the negated caption generation relies on the other captions, we can generate diverse and different negated captions in every training epoch. We also propose a negation text-to-image generation benchmark, NEG-TTOI, to evaluate the capability of models to avoid generating undesired objects given negated prompts. In this task, a compositional negated caption is given which contains the desired objects and undesired objects, e.g., "A women not holding a dog in the car". The generative model needs to explicitly recognize what needs to be generated and what should be avoided. We show that our proposed data generation and training pipeline can directly benefit the downstream task of text-to-image generation. Our contributions include:

- We propose a novel and efficient trainingtime negation generation pipeline, *TNG-CLIP*, to improve CLIP's negation understanding by generating dynamic and diverse negation samples during training without the need for LLMs and pre-defined negation datasets.
- We propose the first benchmark for negationaware text-to-image generation task, NEG-TTOI, which contains diverse and abundant samples to evaluate model's negation understanding capability.
- We offer extensive experiments to demonstrate that *TNG-CLIP* achieves SOTA performance on diverse negation-aware downstream tasks including image-to-text matching, textto-image retrieval, and image generation, indicating its robustness across these tasks.

2 Related Works

While recent foundation models, including LLMs and VLMs, have achieved remarkable success across diverse downstream tasks, their ability to handle negation semantics remains limited. In the scope of large-scale foundation models, the study of negation understanding starts from languageonly setting, where large language models, instead of vision-language models, are focused. Truong et al. shows the LLM's insensitivity of negation by evaluating SOTA LLMs (Brown et al., 2020; Ouyang et al., 2022; Chung et al., 2022) on diverse text-only negation benchmarks (Hossain et al., 2020; Geiger et al., 2020; Truong et al., 2022). Zhang et al. mentions that scaling-up the size of LLM fails to tackle negation tasks. Also, Varshney et al. analyze and tackle the issue of negation in LLM hallucinations, which also emphasizes the significance of negation understanding in LLMs.

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

On the other hand, the negation study in VLMs is mainly focused on CLIP (Radford et al., 2021). For example, Quantmeyer et al. conduct experiments and visualize where and how does CLIP model process negation information in each layers. To make CLIP model understand negation, methods (Park et al., 2025; Singh et al., 2024; Alhamoud et al., 2025a) adopt LLMs to generate negation caption, based on existing image-text pair datasets, to fine-tune the CLIP for negation understanding. However, generating million-scale negation caption with LLM is extremely time- and computeconsuming, and the negation caption is associated with fixed negation object. For example, when a image is paired with the negation caption "A dog not with a boy", the word "boy" can be substituted with plenty of potentially-existing objects such as "cat", "ball", "food" and so on.

Instead of relying on a fixed and stationary dataset throughout the training process, some methods explore the application of dynamic and nonstationary datasets during training process (Wang et al., 2019; Jiang et al., 2024; Böther et al., 2025; Cheng et al., 2025), which is an effective strategy to improve model robustness, generalization, and training efficiency. Inspired by the idea of dynamic dataset training, we generate similar but different negation captions for the same image in every epoch of training, which enhances the diversity of the dataset. Thus, models can learn negation semantics via the absence of multiple negation objects to improve robustness and generalization.



Figure 2: **Training Procedure of** *TNG-CLIP*. The diagram shows the data generation pipeline during the training for one sample in the batch. For an image-text pair, P_o , the most similar image pair, P_s is selected by the cosine similarity of their embedded image features. The captions from P_o and P_s are used to find the negation object and generate two types of negation captions. The final image-text set, S_i , for i^{th} image-text pair will be composed of one image, I_i , one original caption, T_{o_i} , one compositional negation caption, T_{nc_i} , and one full negation caption, T_{nf_i} from another random sample.

3 Training-Time Negation Data Generation for Negation Understanding

170

171

172

174

175

178

179

182

183

184

To make CLIP learn negation semantics with diverse datasets and without the burden of time- and compute-consuming LLM-based negation caption generation, we present our novel training pipeline, Training-Time Negation Data Generation for CLIP (*TNG-CLIP*), such that we generate image-text sets with form $<I, T_o, T_{nc}, T_{nf}>$, from the given image-text pair $<I, T_o>$, where I and T_o represent the provided image and the original (non-negation) caption in the image-text pair dataset, while T_{nc} and T_{nf} represent the two types of generated negation captions: compositional negation caption and full negation caption, discussed in Sec 3.1.3.

3.1 Training time data augmentation

We propose a novel negation data-generation pipeline that the negation captions are formed during each batch of training procedure. The negation data generation pipeline for one image in the batch is shown in Figure 2. Overall, for a given image-text pair, P_o , we will first find another similar image-text pair, P_s , select the negation object, O_n , and generate corresponding negation captions, T_{nc} and T_{nf} with the randomly-chosen negation pattern template and form the image-text set, S.

193

194

195

197

198

200

201

202

203

204

206

207

208

209

210

211

212

213

214

215

216

3.1.1 Find similar image-text pairs

To form a semantically reasonable compositional negation caption, T_{nc} , we need to find a proper negation object, O_n , that can be potentially fitted into the original caption, T_o . For example, we want T_{nc} to be "A dog running with no boy around", instead of "A dog running with no whale around", which is semantically unlikely. Previous methods (Park et al., 2025; Alhamoud et al., 2025b) acquire the proper negation object, O_n , through the reasoning of LLM to find the possible object that might appear in the image but is actually absent. For efficiency, we avoid the use of an LLM, and propose to find the possible O_n of the image-text pair, P_o , from its most similar images-text pair, P_s , in the same batch. Thus, the first step is to find the P_s for every P_o via cosine similarity, between the embedded image features.

Given a visual encoder $E_v(\cdot)$, a batch of images I_b is encoded into the corresponding visual features

$$V_b = E_v(I_b), V_b \in \mathbb{R}^{B \times D}, \tag{1}$$

265 266

267

269

270

271

272

273

274

275

276

277

278

280

281

283

284

285

287

288

290

291

293

294

297

298

299

300

301

302

303

304

306

264

where B is the batch size and D is the hidden dimension of image feature. For i^{th} image feature, V_{b_i} , we apply cosine similarity

$$V_{bs_i} = \arg\max_{V_i} cos_sim(V_{b_i}, V_j)$$
(2)

to find the most similar image feature, V_{bs_i} , and keep track of the most similar image-text pair, P_{s_i} , associated with the image feature V_{bs_i} .

3.1.2 Select negation object

217

218

219

226

232

233

241

242

247

251

253

254

257

258

259

260

263

After having P_s for each image-text pair, P_o , we aim to find the negation object, O_n , exists in P_s 's caption that does not exist in the caption of P_o . For caption in P_s , we employ Natural Language Tool Kit (Bird et al., 2009) to extract the POS tag of every word, and only keep those represent nouns. To avoid selecting the object which is too semantically close to the words in original caption and cause conflict, we use WordNet (Miller, 1995) and its hand-curated symbolic network to select the negation object, O_n , with furthest semantics to those words in the original caption.

3.1.3 Template-based negation caption generation

For every T_o and O_n , we employ randomly-chosen negation templates to generate two different types of negation captions: **compositional negation caption**, T_{nc} , and **full negation caption**, T_{nf} . While the compositional negation caption helps model align image with partial negation of a relevant caption, full negation caption makes the image align with the negation semantics of an unrelated caption.

1. Compositional Negation Caption: The negation caption is in the format of "A < negation > B", where A denotes the original caption, T_o , B denotes the negation object, O_n , and < negation > represents the negation template that combines the two. For example, let A denotes "A dog playing a ball.", B denotes "Boy", and < negation > denotes "There is {caption}, but not a {obj} around." The final compositional negation caption, T_{nc} , is "There is a dog playing a ball, but not a boy around." To make the generated captions diverse, we use GPT-40 (OpenAI et al., 2024) to generate 46 different negation patterns.

2. Full Negation Caption: The negation caption is in the format of *<negation> A*, which is the negation of the entire caption. We use GPT-40 to generate 18 different negation pattern. All the negation patterns and the prompt for GPT-40 to generate them are attached in Appendix A.5.

3.1.4 Form new image-text set

Given the original caption, T_o , compositional negation caption, T_{nc} , and full negation caption, T_{nf} , we can now construct the final image-text set, S, for training. For each image I_i , we associate it with the original caption, T_{o_i} , the compositional negation caption, T_{nc_i} , and the full negation caption, T_{nf_j} , $j \neq i$. Please note that we randomly pick the full negation caption, T_{nf_j} , from other imagetext pairs, P_j . This is because we want to align the negation of the irrelevant captions to the image and contrast the negation of the relevant caption. Finally, the image-text set, S, is denoated as

$$\mathrm{Image}_i \leftrightarrow \left\{ \begin{array}{ll} \mathrm{Original}_i \\ \mathrm{Compositional\ Negation}_i \\ \mathrm{full\ negation}_j, j \neq i \end{array} \right. 2$$

3.2 Asymmetric noise-augmented objective

After negation image-text set generation, each image is associated with three captions, which makes the image-text pair imbalanced. Thus, the image-totext loss, \mathcal{L}_{i2t} , and text-to-image loss, \mathcal{L}_{t2i} , become asymmetric. We redefine the functionality of both unidirectional loss to serve different purpose.

Text-to-Image Objective Given that we have three captions for one image, the similarity matrix will be in shape of $3N \times N$, where N denotes the number of the images. We calculate the \mathcal{L}_{t2i} in a single objective by applying same image alignment to the three captions. The text-to-image objective function is defined as:

$$\mathcal{L}_{t2i} = -\frac{1}{3N} \sum_{j=0}^{3N-1} \log \left(\frac{\exp\left(S_{j, \lfloor \frac{j}{3} \rfloor} / \tau\right)}{\sum_{i=0}^{N-1} \exp\left(S_{j, i} / \tau\right)} \right),$$

where $S_{j,i}$ denotes the similarity between caption j and image i.

Image-to-Text Objective Aligning each image with a negation caption, specifically negation object, is out-of-distribution for pre-trained CLIP because CLIP, which has seen only image-text pairs in which almost all textual components are visually grounded, with no explicit representation of negation. As a result, the pre-trained model struggles to align negation semantics or irrelevant objects with the image. Fine-tuning pre-trained model on such OOD task might lead to worse performance,

397

352

353

because fine-tuning can achieve worse accuracy, by 307 overfitting, when the pretrained models are good 308 and the downstream task distribution shift is large, supported by theory from (Kumar et al., 2022). To solve the above obstacle of overfitting, we introduce label noise to improve the generalization and 312 robustness of the model, inspired by the related 313 works (Rolnick et al., 2018; Xie et al., 2020; Chen 314 et al., 2025). We modified the image-to-text loss 315 such that the text labels are randomly aligned with 316 the image to introduce noise to the objective func-317 tion. The \mathcal{L}_{i2t} is: 318

$$\mathcal{L}_{i2t} = -\frac{1}{N} \sum_{i=0}^{N-1} \log \left(\frac{\exp(S_{i,y_i}/\tau)}{\sum_{j=0}^{3N-1} \exp(S_{i,j}/\tau)} \right),$$

where $y_i \sim \mathcal{U}(\{0, 1, \dots, 3N - 1\})$ is a random selected label across all the captions labels.

319

324

328

334

336

337

Combined Objective By introducing noise to \mathcal{L}_{i2t} , we only have uni-directional \mathcal{L}_{t2i} helping align negation captions to image. This approach is possible because we freeze the visual encoder during the training, following previous works (Singh et al., 2024; Park et al., 2025). Because the visual encoder is fixed, the visual feature is not updated during image-to-text alignment training, and the model only learn to update text features closer to the pre-trained visual features. The final objective function is then defined as:

$$\mathcal{L} = \frac{1}{2} (\mathcal{L}_{i2t} + \mathcal{L}_{t2i}).$$

The further analysis of the objective function is presented in Appendix A.1.

4 Negation Text-to-Image Generation Benchmark

While negation is an essential part of natural language understanding, a well-designed image gen-339 erative model should be capable of understanding 340 what to generate and what to avoid. To analyze 341 the generative models' performance on negation prompts, Park et al. proposed negation-aware image generation experiments with only 107 negation prompts, containing simple naive negation pattern 345 of "no", "not", "without". To enable systematic 347 analysis, we design the first negation-based textto-image generation benchmark, NEG-TTOI, with examples in Table 11. It contains 2000 evaluation samples in the form of $\langle p, q_p, q_n, a_p, a_n \rangle$, where p is the prompt mentioning both desired and undesired 351

objects, q_p is positive question about the existence of desired objects, q_n is the negative question about the absence of undesired objects, and a_p and a_n are the answer to q_p and q_n .

4.1 Negation prompts generation pipeline

We follow the procedure of previous works (Park et al., 2025; Alhamoud et al., 2025a) to generate prompts and questions via LLM. We use LLM instead of our negation generation pipeline in Sec 3 because (1) the scale of our evaluation benchmark is much smaller than the scale of training dataset, and (2) we only generate the benchmark prompts and questions once, without the necessity of iterative negation data generation over epochs, which makes the LLM time- and compute-affordable.

We use the MS-COCO Caption (Chen et al., 2015) as the base dataset. The goal of our caption generation pipeline is to transform each caption, which describes the existing scene or objects in the image, into a negation-style caption in which certain elements are explicitly described as absent. To efficiently manipulate the caption with complicated semantics, we leverage GPT-40 (OpenAI et al., 2024) in a multi-step manner from negation prompt generation, evaluation questions generation and quality verification.

- 1. Negation Prompt Generation: For every input caption, we ask LLM to identify a random scene or object that is mentioned in the original caption. The selected scene or object will be used as the negation object to generate negation caption. Once we have the original caption and the negation object, we prompt LLM to rewrite the original caption such that the object should be semantically absent from the original caption.
- 2. Evaluation Question Generation For every negation prompt, we prompt LLM to identify the positive semantics and negative semantics in the sentence while discard the negation pattern. For example, given a negation caption "A dog playing a yellow ball while there is no man walking around", the positive semantics will be "A dog playing a yellow ball", while the negative semantics will be "man walking around". Both the positive semantics and negative semantics are combined with "Is there...?" to form the questions q_p and q_n.
- **3. Question Quality Verification** Although GPT-40 is one of the SOTA LLMs for semantic under-

Model	Avg.	Affirmation	Negation	Hybrid	R@5	Neg-R@5
CLIP (Pretrained)	16.28	21.89	16.89	9.99	54.76	47.92
CoN-CLIP	15.70	0.05	36.73	11.97	51.91	48.22
NegCLIP	10.21	9.97	19.76	1.83	68.73	64.41
CLIP (CC12MNegFull)	46.9	56.49	41.71	42.29	54.20	51.90
TNG-CLIP (Ours)	52.5	68.75	44.75	43.29	62.00	61.11

Table 1: Result on Negbench MSCOCO image dataset on image-to-text matching and text-to-image retrieval tasks. **R@5** refers to the Top-5 accuracy on original (non-negation) MSCOCO-Caption dataset, while **Neg-R@5** refers to the Top-5 accuracy on negation MSCOCO-Caption dataset from NegBench.

Model	Avg.	Affirmation	Negation	Hybrid
CLIP (Pretrained)	14.47	31.96	8.34	14.97
CoN-CLIP	22.36	0.01	27.67	24.14
NegCLIP	8.50	22.58	8.62	4.08
CLIP (CC12MNegFull)	52.65	73.75	35.69	62.34
TNG-CLIP (Ours)	59.23	85.92	36.39	72.80

Table 2: Result of Negbench image-to-text matching on VOC2007 image dataset

standing, it still might generate text that are se-401 402 mantically incorrect. Thus, verification is necessary to prevent the improper generation. Given 403 the negation prompt, p, positive question, q_p , 404 and negative question q_n , we prompt the LLM 405 to ask whether the semantics in the q_p is stated 406 positively in p, and whether the semantics in the 407 q_n is stated negatively in p with the negation 408 semantics. If the LLM's answer for both ques-409 tions are correct, the negation data sample will 410 be kept, otherwise it will be discarded. 411

In the end, NEG-TTOI contains 2000 valid samples, selected from 2500 candidates.

4.2 Evaluation metrics

412

413

414

427

428

429

415 Unlike image-text matching or retrieval tasks such that the explicit ground truth can be found, evaluat-416 ing image generation task is relatively subjective. 417 Inspired by (Park et al., 2025; Hu et al., 2023), we 418 employ GPT-40 (OpenAI et al., 2024) to evaluate 419 the existence and absence of the objects. Given a 420 image generated using negation caption as prompt 421 and the positive question and negative question, we 422 evaluate the model's generation quality via the met-423 ric of Compositional Accuracy: it's True if the 424 LLM answers "yes" on positive question and "no" 425 on negative question at the same time. 426

5 Experiments

To show the capability of our proposed method on multiple downstream tasks, we evaluate our model

on negation tasks including image-to-text match-430 ing, text-to-image retrieval and text-to-image gen-431 eration. Our goal is to assess TNG-CLIP's negation 432 semantics understanding via multiple benchmarks 433 and show its generalization and capacity on diverse 434 negation-based scenarios. In the paper, all experi-435 ments are performed on a single Nvidia A40 GPU 436 with batch size of 128 and learning rate of 5e-6. 437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

5.1 Matching & retrieval evaluation

To evaluate the negation understanding ability of *TNG-CLIP*, we present the experiments on image-to-text matching and text-to-image retrieval tasks.

Benchmarks We employ the following benchmarks to evaluate the model's performance:

- Valse-Existence (Parcalabescu et al., 2022) benchmark evaluates the model's performance on negation imaget-to-text matching task. Given a image and two text description about the presence and absence of an object in the image, *e.g. "There is animal in the image"/"There is no animal in the image"*, the model should select the best-matched text.
- NegBench (Alhamoud et al., 2025b) benchmark is a comprehensive benchmark to evaluate the negation understanding of models on variant image-to-text matching and text-to-image retrieval tasks. It includes negation-based matching tasks based on both MS-COCO(Chen et al., 2015) and

VOC2007(Everingham et al.) datasets, a textto-image retrieval task based on MS-COCO evaluation dataset, where the captions are converted into compositional negation style. In the matching task, images are paired with four different captions of three categories: *Affirmation* for *"It include A and B."*, *Negation* for *"Does not include A and B."*, and *Hybrid* for *"Include A but not B."*.

459

460

461

462

463

464

465

466

467

Model	Accuracy
CLIP (Pretrained)	65.16
NegCLIP	73.22
CoN-CLIP	74.15
CLIP (CC12MNegFull)	76.21
NegationCLIP	80.15
TNG-CLIP (Ours)	81.64

Table 3: Valse-Existence Image-to-Text Matching

Baselines To evaluate the performance of our 468 method, we compare it against several existing 469 470 baseline methods for CLIP's negation understanding, including pretrained-CLIP (Radford et al., 471 2021), NegCLIP (Yuksekgonul et al., 2023), CoN-472 CLIP (Singh et al., 2024), and CLIP fine-tuned on 473 CC12M-NegFull (Alhamoud et al., 2025a). For fair 474 comparison, all of the methods are initialized based 475 on pre-trained CLIP ViT-B/32 model. 476

Comparison Experiments We present the 477 matching and retrieval task of NegBench-478 MSCOCO in table 1 and the matching task of 479 480 NegBench-VOC2007 in table 2. From the tables, we observe that previous methods are lack of gener-481 alization on negation-based tasks, but only focus on 482 the negation understanding of specific tasks. For ex-483 ample, CoN-CLIP's performance on matching (af-484 firmation) task is 0.05 and 0.29 on MSCOCO and 485 VOC2007 datasets, which indicates that the method 486 is biased such that it sacrifices the CLIP's perfor-487 mance on non-negation performance for negation 488 improvement. For NegCLIP, even though it get 489 the best score on retrieval task, we observe that the 490 affirmation performance is lower than that of the 491 pretrained-CLIP, and its performance on matching 492 493 (hybrid) is low. On the other hand, the CC12M-NegFull fine-tuned CLIP presents general improve-494 ment of different tasks, indicating its capability of 495 diverse negation tasks. Our method, TNG-CLIP, 496 even though slightly underperforms the NegCLIP 497

Strategy	Avg. Acc.	
dynamic dataset	51.61 ± 0.96	
fixed dataset	49.52 ± 1.27	

Table 4: Effect of using dynamic dataset. Evaluation on NegBench-MSCOCO image-to-text matching task.

model on retrieval tasks, achieves SOTA performance on all the matching tasks, shows its generalization and high-performance on diverse scenarios. 498

499

500

501

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

Similarly, the evaluation on Valse-Existence dataset, in Table 3, further proofs *TNG-CLIP*'s, capability of negation understanding. While the benchmark is first used by *NegationCLIP* (Park et al., 2025) and achieves promising result of 80.15 on CLIP ViT-B/32 based models, our method gets better performance, 81.64, which is higher than all other negation-understanding CLIP baselines.

Effectiveness of Dynamic Dataset The trainingtime data generation pipeline generates the negation caption based on the other image-text pairs in the same batch, which makes the negation caption of same image different in every epoch. We analyze the effect of such dynamic dataset and compare how the performance differs from using fixed dataset. We store the image-text set, S, generated in each training epoch for every epoch as the fixed dataset. We then use the fixed dataset to replace the data generation pipeline to fine-tune the CLIP model. To get statistically significant comparison result, we repeat the TNG-CLIP's training procedure for 10 times and use 10 fixed dataset collected from different training epochs to finetune pre-trained CLIP with same objective function and hyper-parameters. We present the mean and standard deviation in Table 4. We observe that the performance of TNG-CLIP is higher than using fixed dataset, and the standard deviation is also smaller than the fixed one. We explain such phenomenon as the CLIP's fine-tuning on fixed dataset constrains the model's negation understanding to specific *<caption*, *negation object>* pair, thus harms the generalization of the model on negation tasks, leading to lower mean accuracy. At the same time, the data variance among every epoch for TNG-CLIP works as a natural regularization to prevent overfitting and memorizing incorrect correlation, thus lead to smaller standard variance.

More analytic experiments are in Appendix A.4 and A.3.

Model	Arch.	Acc.
SD-1.5	ViT-L/14	32.60
SDXL-1.0	ViT-L/14	27.45
SD-1.5 w/ CoN-CLIP	ViT-L/14	28.40
SD-1.5 w/ TNG-CLIP (ours)	ViT-L/14	45.65
pretrained-CLIP + proj	ViT-B/32	28.25
NegCLIP + proj	ViT-B/32	33.85
CoN-CLIP + proj	ViT-B/32	24.05
CC12MNegFull + proj	ViT-B/32	36.95
TNG-CLIP + proj (ours)	ViT-B/32	41.70

Table 5: Image Generation on NEG-TTOI benchmark

544 545 546 547 548 549 550 551 552

553

555

557

560

561

564

566

570

572

574

542

543

5.2.1 CLIP for Image Generation Task

Although CLIP model is mostly used to do imagetext matching tasks, it can be applied to text-toimage generation tasks indirectly. For example, the text encoder from stable diffusion model is the original copy of CLIP ViT-L/14's text encoder (Rombach et al., 2022). To evaluate the negation understanding of CLIP in text-to-image generation field, Park et al. provides a simple yet effective way, by replacing the original text encoder from stable diffusion model with their proposed negation-aware CLIP. This direct substitution is possible because they fine-tune only the text encoder, preserving the original image embedding space and maintaining the text feature alignment with it.

5.2.2 Experiment Setup

Following the strategy mentioned above, we finetuned our *TNG-CLIP* from pretrained CLIP ViT-L/14 model, and replace the original stable diffusion model's text encoder with ours.

However, most baseline methods are fine-tuned only on CLIP ViT-B/32 model, it is difficult to do the direct substitution due to the mismatch of output feature dimension. To tackle such issue, we attach a MLP projector after the frozen text encoder, and perform knowledge distillation between CLIP ViT-L/14's text encoder acts and CLIP ViT-B/32's text encoder with projector, to align the output of projected CLIP ViT-B/32 text encoder similar to that of CLIP ViT-L/14 text encoder. We perform add MLP to all the baseline methods and fine-tune the MLP, with text encoder frozen, on the same dataset, MS-COCO Caption (Chen et al., 2015).

5.2.3 Experiment Analysis

The comparison results on NEG-TTOI benchmark are presented in Table 5. The upper table shows the comparison with CLIP ViT-L/14's text encoder architecture. We choose SD-1.5 (Rombach et al., 2022) as the generative model backbone and replace its text encoder with that of ours and *CoN*-*CLIP*'s. All the experiment here are the zero-shot performance on NEG-TTOI benchmark. We observe that among the all, using our *TNG-CLIP*'s text encoder achieves the best accuracy, indicating its outstanding capability of handling negation feature for image generation. On the other hand, the accuracy of *CoN-CLIP* is lower than original stable diffusion model, which shows its deficiency on image generation task. 578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

The lower table presents the accuracy of SD-1.5 by replacing its text encoder with the combination of CLIP ViT-B/32 based architecture and the fine-tuned MLP projector. Noticing that the accuracy of our method using CLIP ViT-B/32's text encoder is 41.70, while that for using CLIP ViT-L/14's text encoder is 45.65, showing that the projected ViT-B/32 text encoder is not as effective as ViT-L/14's text encoder, and is only used for the purpose of providing accessible and fair comparison between the baselines on image generation task. Among the all, our method's text encoder still achieves the best accuracy, and the clip fine-tuned with CC12MNegFull (Alhamoud et al., 2025a) is the second best, similar with its performance in image-text matching tasks.

We provide more detailed image generation task analysis in Appendix A.2.

6 Discussion & Conclusion

In this paper, we focus on the critical problem of improving negation understanding for CLIP. Instead of using pre-generated fixed negation dataset, we propose a training-time negation data generation pipeline to generate dynamic negation caption during the training time, addressing the time- and compute- inefficiency problem of previous dataset. We also show that using dynamic negation caption during the training can improve mdoel's generalization and boost the performance of negation fine-tuned CLIP. On the other hand, we propose the first negation-aware text-to-image generation evaluation benchmark to expand the horizon of negationrelated benchmarks. Overall, our work underscores the negation understanding in the study of vision language model, and call for the wider exploration of negation-aware model in diverse tasks.

7 Limitations

627

637

641

653

664

665

667

670

671

672

673

675

In this paper, we propose a negation-aware CLIP, *TNG-CLIP*, trained via the novel efficient trainingtime negation data generation pipeline. We also propose a negation text-to-image generation benchmark, NEG-TTOI, to evaluate the capability of generative model's performance with negation semantics. However, although we have shown the performance and generalization of *TNG-CLIP* via multiple benchmarks, we see the limit of our paper:

- In the paper, we mainly focus on the negation understanding of CLIP model. As the lack of negation understanding is an overall challenge among all vision language models, further exploration on negation-awareness of diverse VLMs is necessary.
- The training-time negation data generation pipeline is currently limited to image-text pair dataset, which is adopted to apply contrastive learning. Our negation data generation pipeline has the potential to be extended beyond image-text pairs, eg. visual question answering dataset, thus supports the negation-awareness training with objective function other than contrastive loss.

References

- Kumail Alhamoud, Shaden Alshammari, Yonglong Tian, Guohao Li, Philip Torr, Yoon Kim, and Marzyeh Ghassemi. 2025a. Vision-language models do not understand negation. *Preprint*, arXiv:2501.09425.
- Kumail Alhamoud, Shaden Alshammari, Yonglong Tian, Guohao Li, Philip Torr, Yoon Kim, and Marzyeh Ghassemi. 2025b. Vision-language models do not understand negation. *Preprint*, arXiv:2501.09425.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. Natural language processing with Python: analyzing text with the natural language toolkit. " O'Reilly Media, Inc.".
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. *Preprint*, arXiv:2005.14165.

Maximilian Böther, Ties Robroek, Viktor Gsteiger, Robin Holzinger, Xianzhe Ma, Pınar Tözün, and Ana Klimovic. 2025. Modyn: Data-centric machine learning pipeline orchestration. *Proceedings of the ACM on Management of Data*, 3(1):1–30.

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

- Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. 2021. Conceptual 12m: Pushing webscale image-text pre-training to recognize long-tail visual concepts. *Preprint*, arXiv:2102.08981.
- Hao Chen, Zihan Wang, Ran Tao, Hongxin Wei, Xing Xie, Masashi Sugiyama, Bhiksha Raj, and Jindong Wang. 2025. Impact of noisy supervision in foundation model learning. *Preprint*, arXiv:2403.06869.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollar, and C. Lawrence Zitnick. 2015. Microsoft coco captions: Data collection and evaluation server. *Preprint*, arXiv:1504.00325.
- Ziheng Cheng, Zhong Li, and Jiang Bian. 2025. Dataefficient training by evolved sampling.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret
Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi
Wang, Mostafa Dehghani, Siddhartha Brahma, Albert
Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac
Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex
Castro-Ros, Marie Pellat, Kevin Robinson, and 16
others. 2022. Scaling instruction-finetuned language
models. *Preprint*, arXiv:2210.11416.696
697
- Dumitru, Ian Goodfellow, Will Cukierski, and
Yoshua Bengio. 2013. Challenges in represen-
tation learning: Facial expression recognition704705
challenge. https://kaggle.com/competitions/
challenges-in-representation-learning-facial-expression-
Kaggle.707
- M. Everingham, L. Van Gool, C. K. I. Williams,
 J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007
 (VOC2007) Results. http://www.pascalnetwork.org/challenges/VOC/voc2007/workshop/index.html.
 714
- Atticus Geiger, Kyle Richardson, and Christopher Potts. 2020. Neural natural language inference models partially embed theories of lexical entailment and negation. In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 163–173, Online. Association for Computational Linguistics.
- Micah Hodosh, Peter Young, and Julia Hockenmaier. 2013. Framing image description as a ranking task: data, models and evaluation metrics. *J. Artif. Int. Res.*, 47(1):853–899.
- Md Mosharaf Hossain, Venelin Kovatchev, Pranoy Dutta, Tiffany Kao, Elizabeth Wei, and Eduardo Blanco. 2020. An analysis of natural language inference benchmarks through the lens of negation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*,

Ramesh, Gabriel 733 tional Linguistics. try, Amanda As 734 Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Gretchen Krueg Mari Ostendorf, Ranjay Krishna, and Noah A Smith. 735 ing transferable Tifa: Accurate and interpretable text-to-2023. supervision. Pre image faithfulness evaluation with question answer-737 ing. Preprint, arXiv:2303.11897. David Rolnick, An Shavit. 2018. Yiding Jiang, Allan Zhou, Zhili Feng, Sadhika Malladi, label noise. Pre 740 and J Zico Kolter. 2024. Adaptive data optimization: 741 Dynamic sample selection with scaling laws. arXiv Robin Rombach, A preprint arXiv:2410.11820. 742 Patrick Esser, resolution image Alex Krizhevsky. 2009. Learning multiple layers of 743 els. Preprint, ar 744 features from tiny images. Jaisidh Singh, Ishaa 745 Ananya Kumar, Aditi Raghunathan, Robbie Jones, Singh, and Apar Tengyu Ma, and Percy Liang. 2022. Fine-tuning 746 "yes" better: Im can distort pretrained features and underperform out-747 negations. Prep. of-distribution. Preprint, arXiv:2202.10054. 748 Thinh Hung Truon 749 George A. Miller. 1995. Wordnet: a lexical database for english. Commun. ACM, 38(11):39-41. and Trevor Coh naysayers: An a OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, tion benchmarks Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec 753 Thinh Hung Truon Radford, Aleksander Mądry, Alex Baker-Whitcomb, win, Trevor Coh Alex Beutel, Alex Borzunov, Alex Carney, Alex 2022. Not anot Chow, Alex Kirillov, and 401 others. 2024. Gpt-40 756 NLI test suite fo 757 system card. Preprint, arXiv:2410.21276. ings of the 2nd (ter of the Associ 758 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carand the 12th Inte 759 roll L. Wainwright, Pamela Mishkin, Chong Zhang, ral Language P Sandhini Agarwal, Katarina Slama, Alex Ray, John pages 883-894, Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, tational Linguis Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Neeraj Varshney, S Training language models to follow instructions with neet Chatterjee human feedback. Preprint, arXiv:2203.02155. Chitta Baral. 20 hallucinations of Letitia Parcalabescu, Michele Cafagna, Lilitta Murad-Preprint, arXiv: jan, Anette Frank, Iacer Calixto, and Albert Gatt. 2022. Valse: A task-independent benchmark for Rui Wang, Masa vision and language models centered on linguistic 2019. Dynam phenomena. In Proceedings of the 60th Annual Meettraining of neu ing of the Association for Computational Linguistics arXiv:1805.001 (Volume 1: Long Papers), page 8253-8280. Association for Computational Linguistics. 773 Jianxiong Xiao, Ja Oliva, and Anto 774 Junsung Park, Jungbeom Lee, Jongvoon Song, Sang-Large-scale scer won Yu, Dahuin Jung, and Sungroh Yoon. 2025. 775 2010 IEEE Con Know "no" better: A data-driven approach for 776 puter Vision and enhancing negation awareness in clip. Preprint, 777 3492. arXiv:2501.10913. Qizhe Xie, Minh-Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Quoc V. Le. 20 Juan C. Caicedo, Julia Hockenmaier, and Svetlana dent improves Lazebnik. 2016. Flickr30k entities: Collecting 781 arXiv:1911.042 region-to-phrase correspondences for richer imageto-sentence models. Preprint, arXiv:1505.04870. Mert Yuksekgonul, Vincent Quantmeyer, Pablo Mosteiro, and Albert Gatt. Dan Jurafsky, ar 2024. How and where does clip process negation? vision-language and what to do a 786 Preprint, arXiv:2407.10488. 10

pages 9106–9118, Online. Association for Computa-

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	787
Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	788
try, Amanda Askell, Pamela Mishkin, Jack Clark,	789
Gretchen Krueger, and Ilya Sutskever. 2021. Learn-	790
ing transferable visual models from natural language	791
supervision. <i>Preprint</i> , arXiv:2103.00020.	792
David Rolnick, Andreas Veit, Serge Belongie, and Nir	793
Shavit. 2018. Deep learning is robust to massive	794
label noise. <i>Preprint</i> , arXiv:1705.10694.	795
Robin Rombach, Andreas Blattmann, Dominik Lorenz,	796
Patrick Esser, and Björn Ommer. 2022. High-	797
resolution image synthesis with latent diffusion mod-	798
els. <i>Preprint</i> , arXiv:2112.10752.	799
Jaisidh Singh, Ishaan Shrivastava, Mayank Vatsa, Richa	800
Singh, and Aparna Bharati. 2024. Learn "no" to say	801
"yes" better: Improving vision-language models via	802
negations. <i>Preprint</i> , arXiv:2403.20312.	803
Thinh Hung Truong, Timothy Baldwin, Karin Verspoor,	804
and Trevor Cohn. 2023. Language models are not	805
naysayers: An analysis of language models on nega-	806
tion benchmarks. <i>Preprint</i> , arXiv:2306.08189.	807
Thinh Hung Truong, Yulia Otmakhova, Timothy Baldwin, Trevor Cohn, Jey Han Lau, and Karin Verspoor. 2022. Not another negation benchmark: The NaN-NLI test suite for sub-clausal negation. In <i>Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 883–894, Online only. Association for Computational Linguistics.	808 809 810 811 812 813 814 815 816 816
Neeraj Varshney, Satyam Raj, Venkatesh Mishra, Ag-	818
neet Chatterjee, Ritika Sarkar, Amir Saeidi, and	819
Chitta Baral. 2024. Investigating and addressing	820
hallucinations of Ilms in tasks involving negation.	821
<i>Preprint</i> , arXiv:2406.05494.	822
Rui Wang, Masao Utiyama, and Eiichiro Sumita.	823
2019. Dynamic sentence sampling for efficient	824
training of neural machine translation. <i>Preprint</i> ,	825
arXiv:1805.00178.	826
Jianxiong Xiao, James Hays, Krista A. Ehinger, Aude	827
Oliva, and Antonio Torralba. 2010. Sun database:	828
Large-scale scene recognition from abbey to zoo. In	829
2010 IEEE Computer Society Conference on Com-	830
puter Vision and Pattern Recognition, pages 3485–	831
3492.	831
Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and	833
Quoc V. Le. 2020. Self-training with noisy stu-	834
dent improves imagenet classification. <i>Preprint</i> ,	835
arXiv:1911.04252.	836
Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri,	837
Dan Jurafsky, and James Zou. 2023. When and why	838
vision-language models behave like bags-of-words,	839
and what to do about it? <i>Preprint</i> , arXiv:2210.01936.	840

842 843 844

845

846

852

857

867

871

872

873

875

876

877

880

884

885

841

Yuhui Zhang, Michihiro Yasunaga, Zhengping Zhou, Jeff Z. HaoChen, James Zou, Percy Liang, and Serena Yeung. 2023. Beyond positive scaling: How negation impacts scaling trends of language models. *Preprint*, arXiv:2305.17311.

A Appendix

A.1 Ablation Study of Asymmetric Noise-Augmented Objective

In order to train the negation-aware CLIP for diverse tasks, we propose a novel asymmetric noiseaugmented loss that different from the original contrastive loss of CLIP. We exam the contribution of each component in this novel objective function with analytic ablation study. Within the objective function, we split its component to four parts based on the functionality of each:

- compositional alignment refers to align the compositional negation caption to the image in \mathcal{L}_{t2i} .
- **full alignment** refers to align the full negation caption to the image in the \mathcal{L}_{t2i} .
- original alignment refers to align the original caption to the image in \mathcal{L}_{t2i} .
- **noise alignment** refers to align the randomchose caption to the image in \mathcal{L}_{i2t} .

The analytic results are presented in Table 6, with the evaluation on Negbench-MSCOCO matching and Negbench-MSCOCO Retrieval tasks. From the table, we observe that by eliminating compositional negation, full negation and original caption from \mathcal{L}_{t2i} separately, the corresponding performance in matching task drops. For example, without original caption, the affirmation accuracy drops from 68.75 to 60.18. At the same time, the accuracy of negation retrieval tasks remains similar, indicating that the components in \mathcal{L}_{t2i} are not the primary factors for it.

We then analyze the effect of random noise in \mathcal{L}_{i2t} . Instead of letting image random choose caption, we match the image to its corresponding original, compositional negation and full negation captions as three independent experiments. Additionally, we let images to randomly match one of its corresponding original, compositional negation and full negation caption, and even directly delete the \mathcal{L}_{i2t} loss. Through the experiments, we found that without the random noise, performance of retrieval task drops significantly. This matches the

hypothesis we proposed in Sec 3.2 that negation dataset is an OOD task for pre-trained CLIP, the direct fine-tuning may cause worse performance.

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

Lastly, we propose and examine another loss objective: can we split the generated image-text set, S, to form three image-text pairs for each of original, compositional negation and full negation, and apply normal contrastive loss on the three independently? We implement such objective function and present it at the bottom of the table. We observe that by doing so, the performance of both matching task and retrieval task are sub-optimal. While there is also no noise label added to the objective training, the worse result on using independent losses, again, emphasizes the importance of adding noise when fine-tuning CLIP on negation-related dataset.

A.2 More Analysis of Image Generation Experiment

In the image generation task, we observe the inefficiency of original Stable Diffusion and CoN-CLIP in the NEGTTOI benchmark. But why this happens? To further explore that, we evaluate the performance of models with two analytic metrics: Positive Accuracy and Negative Accuracy. Given a prompt "generate A without B", Positive Accuracy measures if the image contains A, and Negative Accuracy measures if the image doesn't contain B. The result is presented in Table 7. In the table, we can observe that for original Stable Diffusion model, the positive accuracy is higher than that of using our method or CoN-CLIP as text encoder, but the negative accuracy is much lower. This explicitly shows that the original text encoder cannot process negation semantics to help avoid the generation of unwanted objects. On the other hand, adopting CoN-CLIP as text encoder can significantly boost the negative accuracy, but at the same time, its performance on positive accuracy becomes low. This indicates the CoN-CLIP model is a biased model towards negation-understanding, while ignores the generalization on other non-negation tasks.

A.3 Non-Negation Generalization on Image Classification

Although TNG-CLIP is specifically designed for negation understanding, it is important to ensure that its performance on non-negation tasks remains intact, in another word, it should not suffer from catastrophic forgetting on tasks that the original pre-trained CLIP model was capable of handling. Inspired by the experiments from (Singh

Model	Avg.	Affirmation	Negation	Hybrid	Neg-R@5
TNG-CLIP	52.50	68.75	44.75	43.29	61.11
	Abl	ation of Caption C	ategory		
w/o compositional	48.15	65.45	38.02	40.10	56.66
w/o full	51.31	75.63	24.75	51.44	60.79
w/o original	46.66	60.18	45.82	37.79	59.91
Ablation of Noise					
\mathcal{L}_{i2t} : original	52.49	81.93	16.09	56.66	45.32
\mathcal{L}_{i2t} : compositional	50.13	78.05	11.97	57.40	45.39
\mathcal{L}_{i2t} : full	40.29	45.12	44.11	31.85	48.58
\mathcal{L}_{i2t} : random of three	47.92	58.66	43.26	41.10	50.66
w/o \mathcal{L}_{i2t}	46.24	59.54	36.04	42.29	50.49
independent losses	46.49	55.41	43.74	40.05	50.19

Model Positive Negative Arch. SD-1.5 41.95 ViT-L/14 80.85 ViT-L/14 87.05 SDXL-1.0 32.30 SD-1.5 w/ CoN-CLIP ViT-L/14 46.25 72.50 SD-1.5 w/ TNG-CLIP (ours) 75.80 ViT-L/14 63.05 45.65 67.25 pretrained CLIP + proj ViT-B/32 ViT-B/32 67.80 52.10 NegCLIP + proj CoN-CLIP + proj ViT-B/32 39.45 72.65 ViT-B/32 53.76 CC12MNegFull + proj 71.80 TNG-CLIP + proj ViT-B/32 63.65 68.50

Table 6: Ablation Study on NegBench MSCOCO matching task

Table 7: Image Generation on Neg-TtoI benchmark

et al., 2024), we conduct the zero shot image 939 classification on TNG-CLIP and pre-trained CLIP with eight diverse benchmarks: FER2013 (Du-941 mitru et al., 2013), Flickr-8K (Hodosh et al., 942 2013), Flickr-30K (Plummer et al., 2016), MS-943 COCO (Chen et al., 2015), SUN397 (Xiao et al., 2010), VOC2007 (Everingham et al.), CIFAR-945 10 (Krizhevsky, 2009), CIFAR-100 (Krizhevsky, 2009). The top1 and top5 accuracy score is presented in Figure 3. In the figure, we observe that the 948 zero-shot performance of TNG-CLIP remains similar with that of pre-trained CLIP, indicating there is no catastrophic forgetting or overfitting to our 951 952 proposed method. Surprisingly, we also observe that in some cases, such as Flickr-8K, Flickr-30K, MS-COCO and VOC2007 benchmarks, the TNG-954 CLIP outperforms the pre-trained CLIP, illustrating 955 956 improving the negation understanding can improve

model's performance on general tasks.

A.4 Time-Efficiency Test

As we generate data samples during the training stage, does the generation pipeline significantly slower the training process and becomes timeconsuming? We compares the average training time per batch on the same GPU device, Nvidia-A40, with and without the data generation pipeline in Table 8. For every batch, the data generation pipeline takes 0.13 sec, which is only 2.55% slower than without using the data generation pipeline. Thus, adding the data generation pipeline to the training is still time-efficient. 957

958

959

960

961

962

963

964

965

966

967

968

969

970

A.5 Template-based Negation Pattern

During the negation caption generation, we use pre-
defined LLM-generated negation pattern template971to convert the original caption and negation object973



Figure 3: The zero shot image classification accuracy of pre-trained CLIP and TNG-CLIP on eight image classification benchmarks.

Strategy	Time (sec)
w/o data generation	4.97
w/ data generation	5.10
data generation	0.13

Table 8: Time Efficiency for Data Generaiton

974to compositional negation caption and full negation975caption. We present the template we used here in976Table 9 and Table 10.

There's no {cap} in the image.	No {cap} is included in the image.
There is not {cap} in the image.	The image does not have {cap}.
No {cap} is present in the image.	{cap} is not present in the image.
{cap} is absent.	No {cap} is present.
There isn't any {cap}.	Not a single {cap} can be seen.
The image is without {cap}.	The image is lacking {cap}.
There appears to be no {cap} in the image.	The image does not contain {cap}.
There does not exist {cap} in the image.	There is nothing about {cap}.
There isn't any {cap}.	No {cap} is seen in the image.

Table 9: Templates for full negation caption generation, we replace the *cap* with the provided original caption.

{cap} with no {obj}.	{cap} without {obj}
{cap} that do not have {obj}.	{cap} having no {obj}.
{cap} not include {obj}.	{cap} excluding {obj}.
{cap}, but no {obj} are present.	{cap}, though no {obj} can be seen.
{cap} without any {obj} in sight.	{cap} yet no {obj} are nearby.
{cap} but no {obj} are visible.	{cap} and no {obj} are anywhere around.
{cap}, without any {obj} in the vicinity.	{cap}, with no {obj} in the surroundings.
{cap}, but no {obj} are in the area.	{cap}, and no {obj} can be found nearby.
{cap} in the absence of {obj}.	{cap}, where no {obj} are present.
{cap} with an absence of {obj}.	{cap}, as no {obj} are around.
{cap}, while lacking any {obj}.	{cap} but no {obj} are engaging.
{cap} with no {obj} participating.	{cap} yet no {obj} are interacting.
{cap}, as no {obj} are involved.	{cap}, while {obj} remain absent from the
	scene.
{cap} though no {obj} can be spotted.	{cap} where no {obj} are noticeable.
{cap} but no {obj} are detectable.	{cap}, as no {obj} are apparent.
{cap}, with no sight of any {obj}.	No {obj} is visible, but {cap}.
No {obj} can be seen, while {cap} happens.	No {obj} is present, yet {cap} continues.
No {obj} appears in sight, but {cap} unfolds.	Not a single {obj} is noticeable, but {cap}.
No trace of {obj} can be found, while {cap}	No sign of {obj} is apparent, but {cap} is hap-
occurs.	pening.
There is no {obj} in view, but {cap} takes	None of the {obj} are around, yet {cap} con-
place.	tinues.
Not even one {obj} is nearby, but {cap} is	No {obj} exists in the scene, while {cap} hap-
ongoing.	pens.
Absolutely no {obj} is here, yet {cap} re-	Nowhere can {obj} be found, but {cap} is
mains.	evident.
Nowhere in sight is any {obj}, yet {cap} un-	No {obj} is around in the surroundings, but
folds.	{cap} is occurring.

Table 10: Templates for compositional negation caption generation, we replace the cap with the provided original caption and obj with the corresponding negation object.

compositional negation caption	positive question	negative question
A room painted in blue with a white	Is there a room painted in blue	Is there a door?
sink, but no door.	with a white sink?	
A shot inside a kitchen without anyone	Is there a kitchen shown?	Is there anyone present?
present.		
A woman is walking on the sidewalk	Is there a woman walking on the	Is there her dog?
without her dog.	sidewalk?	
A man without a bike at a marina.	Is there a man at a marina?	Is there a bike?
A man is sitting on a bench without a	Is there a man sitting on a bench?	Is there a bicycle nearby?
bicycle nearby.		
There's no kitchen sink beside the door	Is there a door and countertop?	Is there a kitchen sink beside the
and countertop.		door and countertop?
A bathroom without a checkered black	Is there a bathroom?	Is there a checkered black and
and white tile floor.		white tile floor?
A house boat is moored on a riverbank	Is there a house boat moored on	Is there a bike?
with no bikes in sight.	a riverbank?	
A train missing a striped door waiting	Is there a train waiting on a train	Is there a striped door?
on a train track.	track?	
A small airplane flying without a jet	Is there a small airplane flying?	Is there a jet nearby?
nearby.		
A woman is seen without a horse in	Is there a woman in front of a	Is there a horse?
front of a fence with razor wire.	fence with razor wire?	
No vans are traveling over a bridge next	Is there a bridge next to train	Is there a van?
to train tracks.	tracks?	
A person riding a bicycle without any	Is there a person riding a bicycle?	Is there a river nearby?
river nearby.		
No giraffes can be seen in the wood and	Is there a wood and metal fenced	Is there a giraffe?
metal fenced enclosure.	enclosure?	
A row team without a lead woman	Is there a row team?	Is there a lead woman shouting?
shouting.		
A lady is sitting in a room devoid of any	Is there a lady sitting in a room?	Is there a bright pink wall?
bright pink walls.		
A man carrying a plate without any food	Is there a man carrying a plate?	Is there any food on the plate?
on it.		

Table 11: Example from Neg-TtoI negation image generation benchmark