

452 A Limitation, future work, and societal impact

453 A.1 Limitation and future work

454 There are several limitations to this work that future research can further explore. First, we focus
455 our scope on compositionality benchmarks formulated as image-to-text retrieval task. While this is
456 currently the most prevailing evaluation framework, future research can characterize compositionality
457 evaluation as text-to-image retrieval problem, as in the initial efforts considered by [32, 39]. More im-
458 portantly, we hope our work can guide future efforts in creating and ensuring faithful compositionality
459 benchmarks in text-to-image form. Second, in this work, we identify *two* human interpretable dataset
460 biases, the nonsensical and non-fluent biases, which may not cover all dataset artifacts that could
461 possibly be exploited by a model. Future work may utilize more sophisticated techniques to remove
462 spurious dataset artifacts beyond human comprehension [20]. Finally, we focus our evaluations on
463 contrastively learned vision-language models [30]. Future work should include and characterize the
464 compositionality of modern generative vision-language models [1, 5, 21].

465 A.2 Societal impact

466 As vision-language models such as CLIP [30] are becoming the foundation models for many down-
467 stream applications [34, 31], it is imperative to understand the limitations of these models to avoid
468 misuses and undesirable outcomes [6, 2]. Compositionality benchmarks probe a model’s understand-
469 ing of finer-grained concepts, and hence allow us to identify blind spots [42, 45, 26] of seemingly
470 powerful models deemed by standard classification and retrieval benchmarks [9, 23]. Our work fur-
471 ther alleviates common artifacts in existing compositionality benchmarks that result in overestimation
472 of a model’s capability. We hope our proposed benchmark SUGARCREPE leads to more faithful
473 assessment of a vision-language model’s compositionality, and can hence guide more accurate usages
474 of the models. Nevertheless, we note that strong performances on SUGARCREPE do not imply perfect
475 models. We envision SUGARCREPE being one of the many benchmarks used to comprehensively
476 understand the abilities of vision-language models from various aspects.

477 B Implementation details

478 B.1 Hardware information

479 All experiments are run on a machine with an Intel(R) Xeon(R) CPU E5-2678 v3 with a 512G
480 memory and two 48G NVIDIA RTX A6000 GPUs.

481 B.2 Dataset sources

482 We obtain all existing datasets from their original sources released by the authors. We refer readers to
483 these sources for the dataset licenses. To the best of our knowledge, the data we use does not contain
484 personally identifiable information or offensive content.

- 485 • CREPE [26]: We obtain CREPE dataset from its official repository [\[4\]](#).
- 486 • ARO [42]: We obtain ARO dataset from its official repository [\[5\]](#).
- 487 • VL-CheckList [45]: We obtain VL-CheckList dataset from its official repository [\[6\]](#).
- 488 • COCO [23]: We obtain COCO from its official project website [\[7\]](#).

<https://github.com/RAIVNLab/CREPE>

<https://github.com/mertyg/vision-language-models-are-bows>

<https://github.com/om-ai-lab/VL-CheckList>

<https://cocodataset.org/>

489 B.3 Software configuration

490 **Models.** We detail the sources of the pretrained models we use in the paper, and the hyper-parameters
491 used in training our own models.

- 492 • Vera model [24]: We obtain pretrained Vera model released by its author [8].
- 493 • Grammar model [27]: We obtain the Grammar model released by the authors [9].
- 494 • All pretrained CLIP models: We obtain all pretrained CLIP models’ weights from Open-
495 CLIP [10].
- 496 • NEGCLIP [42]: We obtain weights for pretrained NEGCLIP released by the authors [11].
- 497 • Models trained from scratch: We train RN50 based on OpenCLIP codebase and set hyper-
498 rparameters as following: number of warmup steps is 1000, batch size is 256, learning rate
499 is 1e-4, weight decay is 0.1, number of epochs is 30. We augment the original CLIP loss
500 with hard negative captions following NEGCLIP [42].

501 **Evaluations.** We base our evaluation framework on OpenCLIP [16]. We follow all default hyper-
502 parameters used for evaluating models.

503 C Vision-language compositionality benchmarks

504 We provide an overview of existing vision-language compositionality benchmarks below, with Table [7]
505 summarizing the dataset comparisons.

506 C.1 Image-to-text formulation

507 A majority of current benchmarks formulate the evaluation task as image-to-text retrieval problem.
508 These benchmarks generate hard negative texts procedurally through rule-based templates, where
509 each benchmark considers different types of hard negatives.

510 **VL-Checklist [45].** VL-CheckList aims at evaluating vision-language models’ understanding of
511 different objects, attributes, and relationships. It contains REPLACE hard negatives generated by
512 replacing atomic parts of the positive texts with other foils. VL-CheckList further breaks the hard
513 negatives down into more granular categories based on the type of the replaced atomic part, *i.e.*,
514 object, attribute, or relationship.

515 **ARO [42].** ARO focuses on models’ understanding of different relationships, attributes, and order
516 information. It considers SWAP and SHUFFLE hard negatives. SWAP hard negatives are generated by
517 swapping two words in the positive texts; on the other hand, SHUFFLE hard negatives are generated
518 by shuffling words in the positive texts. ARO further divides SWAP hard negatives into attribute or
519 relationship type.

520 **CREPE [26].** CREPE is a large-scale evaluation benchmark that includes three types of hard
521 negatives: REPLACE, SWAP and NEGATE. REPLACE and SWAP hard negatives are generated as
522 in VL-CheckList and ARO. In addition, NEGATE hard negatives are generated by adding negation
523 keywords (*i.e.*, *not* or *no*) to the original positive texts. The hard negatives are not further divided into
524 fine-grained types (object, attribute, or relations).

525 C.2 Text-to-image formulation

526 Complementary to image-to-text formulation, compositionality can as well be evaluated by probing
527 a model to select an image that best matches a given text description, against other hard negative

⁸<https://huggingface.co/liujch1998/vera>

⁹<https://huggingface.co/textattack/distilbert-base-uncased-CoLA>

¹⁰https://github.com/mlfoundations/open_clip

¹¹<https://github.com/mertyg/vision-language-models-are-bows>

Table 7: Summary on vision-language compositionality benchmarks. SUGARCREPE considers image-to-text formulation to enable larger scale evaluation set. In addition, SUGARCREPE considers a wide range of hard negative types. SHUFFLE and NEGATE are omitted as they introduce inevitable biases discussed in Sec. 4.2.

Benchmark	Task Formulation	Scale	Hard Negative Text Type				
			SHUFFLE	REPLACE	SWAP	NEGATE	ADD
VL-CheckList [45]	Image-to-Text	> 1000		✓			
ARO [42]	Image-to-Text	> 1000	✓		✓		
CREPE [26]	Image-to-Text	> 1000		✓	✓	✓	
Winoground [39]	Image-to-Text / Text-to-Image	400			✓		
Cola [32]	Text-to-Image	210			N/A		
SUGARCREPE	Image-to-Text	> 1000		✓	✓		✓

528 images as distractors. Unlike hard negative texts, hard negative images are more difficult to obtain
 529 and thus current text-to-image compositionality benchmarks are smaller at scale.

530 **Winoground [39]**. Winoground is a small dataset manually curated by human annotators. Each
 531 example in the dataset contains two images and two matching captions, where both captions contain
 532 identical words that appear in different orders. Note that Winoground can be used for either image-
 533 to-text or text-to-image retrieval. While the original intention for Winoground is to evaluate vision-
 534 language compositionality, recent work [10] has pointed out that solving the tasks in Winoground
 535 requires not just compositional vision-language understanding, but additionally a suite of other
 536 abilities such as commonsense reasoning, or distinguishing visually difficult images.

537 **Cola [32]**. Cola tests a vision-language model’s ability to select an image that correctly matches a
 538 given caption, against another distractor image with the same objects and attributes but in the wrong
 539 composition. The image pairs are mined from existing datasets. As a result, the final evaluation set is
 540 relatively small in size (210 examples in total).

541 We deem text-to-image evaluation as important as image-to-text evaluation. Future work can explore
 542 approaches to generate or mine compositional hard negative images at scale, as preliminarily explored
 543 in [32, 42].

544 D SUGARCREPE

545 D.1 Taxonomy

546 Figure 6 shows the taxonomy of SUGARCREPE. We first categorize the hard negatives based on
 547 their forms: REPLACE, SWAP, and ADD. We then further divide each type of hard negatives into
 548 finer-grained sub-categories based on the type (object, attribute, or relation) of the atomic concept
 549 altered. SUGARCREPE covers a total of 7 fine-grained hard negative types.

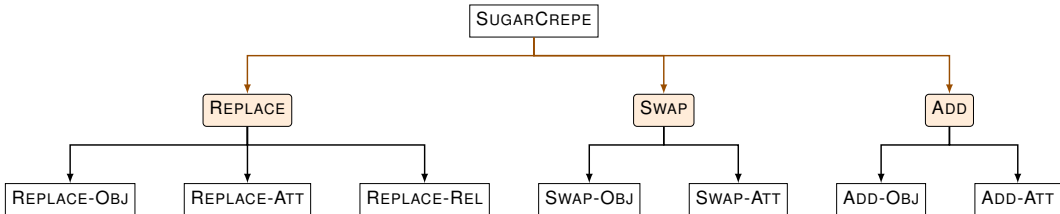


Figure 6: Taxonomy of hard negatives considered in SUGARCREPE.

<p>Given an input sentence describing a scene, your task is to:</p> <ol style="list-style-type: none"> 1. Locate the noun words in the sentence. 2. Randomly pick one noun word. 3. Replace the selected noun word with a new noun word to make a new sentence. <p>The new sentence must meet the following three requirements:</p> <ol style="list-style-type: none"> 1. The new sentence must be describing a scene that is as different as possible from the original scene. 2. The new sentence must be fluent and grammatically correct. 3. The new sentence must make logical sense. <p>Here are some examples:</p> <p>Original sentence: A man is in a kitchen making pizzas. Nouns: ["man", "kitchen", "pizzas"] Selected noun: man New noun: woman New sentence: A woman is in a kitchen making pizzas.</p> <p>Original sentence: a woman seated on wall and birds besides her Nouns: ['woman', 'wall', 'birds'] Selected noun: wall New noun: bench New sentence: A woman seated on a bench and birds besides her.</p>	<p>Given an input sentence describing a scene, your task is to:</p> <ol style="list-style-type: none"> 1. Locate the adjective words describing objects in the sentence. If there are no adjective words, return an empty list. 2. Randomly pick one adjective word. 3. Replace the selected adjective word with a new adjective word to make a new sentence. <p>The new sentence must meet the following three requirements:</p> <ol style="list-style-type: none"> 1. The new sentence must be describing a scene that is as different as possible from the original scene. 2. The new sentence must be fluent and grammatically correct. 3. The new sentence must make logical sense. <p>Here are some examples:</p> <p>Original sentence: a blue bike parked on a side walk. Adjectives: ["blue"] Selected adjective: blue New adjective: red New sentence: a red bike parked on a side walk.</p> <p>Original sentence: The kitchen is clean and ready for us to see. Adjectives: ["clean", "ready"] Selected adjective: clean New adjective: dirty New sentence: The kitchen is dirty and ready for us to see.</p>
--	---

(a) REPLACE-OBJ.

(b) REPLACE-ATT.

<p>Given an input sentence describing a scene, your task is to:</p> <ol style="list-style-type: none"> 1. Find any action or spatial relationships between two objects in the sentence. If there are no such relationships, return an empty list. 2. Randomly pick one relationship. 3. Replace the selected relationship with a new relationship to make a new sentence. <p>The new sentence must meet the following three requirements:</p> <ol style="list-style-type: none"> 1. The new sentence must be describing a scene that is as different as possible from the original scene. 2. The new sentence must be fluent and grammatically correct. 3. The new sentence must make logical sense. <p>Here are some examples:</p> <p>Original sentence: The dining table near the kitchen has a bowl of fruit on it. Relationships: ["near", "on"] Selected relationship: near New relationship: far from New sentence: The dining table far from the kitchen has a bowl of fruit on it.</p> <p>Original sentence: A couple of buckets in a white room. Relationships: ["in"] Selected relationship: in New relationship: outside New sentence: A couple of buckets outside a white room.</p>

(c) REPLACE-REL.

Figure 7: Example prompt templates (black) and outputs (green) from ChatGPT for REPLACE hard negatives.

550 D.2 Hard negative generation procedure and templates

551 To generate hard negatives in SUGARCREPE, we come up with three different prompt templates for
552 the three hard negative types considered: REPLACE, SWAP, and ADD. Each template consists of task
553 instruction for generating the corresponding type of hard negatives and several (7 or more) few-shot
554 demonstrations. We describe the general generation procedure and example prompt templates below
555 and refer readers to our dataset repository for the full prompts used^[12].

556 **Generating REPLACE hard negatives.** To best leverage ChatGPT’s capabilities, we devise a three-
557 step workflow to generate REPLACE hard negatives: (1) We prompt ChatGPT in locating the desired
558 atomic concepts (*e.g.*, objects) in the sentence; (2) We prompt ChatGPT to generate a new concept to
559 replace a randomly selected old concept; (3) We let ChatGPT compose a new sentence by replacing
560 the old concept with the new one. For steps (1) and (3), we prompt ChatGPT with a temperature
561 of 0.0 to get stable outputs. For step (2), however, we diversify the outputs by prompting ChatGPT
562 with a higher temperature of 1.5. With this design, we are able to generate diverse REPLACE hard
563 negatives. Figure 7 shows the example templates and outputs for REPLACE hard negatives.

¹²<https://github.com/RAIVNLab/sugar-crepe>

Given an input sentence describing a scene, your task is to first locate two swappable noun phrases in the sentence, and then swap them to make a new sentence. The new sentence must meet the following three requirements:

1. The new sentence must be describing a different scene from the input sentence.
2. The new sentence must be fluent and grammatically correct.
3. The new sentence must make logical sense.

To complete the task, you should:

1. Answer the question of whether generating such a new sentence is possible using Yes or No.
2. Output the swappable noun phrases.
3. Swap them to make a new sentence.

Here are some examples:

Input: A cat resting on a laptop next to a person.
Is it possible to swap noun phrases in the input sentence to generate a new sentence that is different from the input sentence and makes logical sense? Yes
Swappable noun phrases: laptop, person
Output: A cat resting on a person next to a laptop.

Input: A plate of donuts with a person in the background.
Is it possible to swap noun phrases in the input sentence to generate a new sentence that is different from the input sentence and makes logical sense? Yes
Swappable noun phrases: a plate of donuts, a person
Output: A person with a plate of donuts in the background.

(a) SWAP-OBJ.

Given an input sentence describing a scene, your task is to first locate two swappable adjectives in the sentence describing different objects, and then swap them to make a new sentence. The new sentence must meet the following three requirements:

1. The new sentence must be describing a different scene from the input sentence.
2. The new sentence must be fluent and grammatically correct.
3. The new sentence must make logical sense.

To complete the task, you should:

1. Answer the question of whether generating such a new sentence is possible using Yes or No.
2. Output the swappable adjectives.
3. Swap them to make a new sentence.

Here are some examples:

Input: A girl in a pink shirt holding a blue umbrella.
Is it possible to swap attributes in the input sentence to generate a new sentence that is different from the input sentence and makes logical sense? Yes
Swappable attributes: pink, blue
Output: A girl in a blue shirt holding a pink umbrella.

Input: A girl with a green shirt brushing her teeth with a blue toothbrush.
Is it possible to swap attributes in the input sentence to generate a new sentence that is different from the input sentence and makes logical sense? Yes
Swappable attributes: green, blue
Output: A girl with a blue shirt brushing her teeth with a green toothbrush.

(b) SWAP-ATT.

Figure 8: Example prompt templates (black) and outputs (green) from ChatGPT for SWAP hard negatives.

Given an input sentence describing a scene, your task is:

1. Find the objects in the sentence.
2. Randomly pick one object.
3. Generate a new object that's not in the sentence.
4. Add the new object next to the selected object to make a new sentence.

The new sentence must meet the following three requirements:

1. The new sentence must describe a clearly new and different scene.
2. The new sentence must be fluent and grammatically correct.
3. The new sentence must make logical sense.

Here are some examples:

Original sentence: An elephant standing under the shade of a tree.
Objects: ["elephant", "shade of a tree"]
Selected object: elephant
New object: squirrel
New sentence: An elephant and a squirrel standing under the shade of a tree.

Original sentence: A bench at the beach next to the sea
Objects: ['bench', 'beach', 'sea']
Selected object: bench
New object: umbrella
New sentence: An umbrella and a bench at the beach next to the sea.

(a) ADD-OBJ.

Given an input sentence describing a scene, your task is:

1. Find the objects in the sentence.
2. Randomly pick one object.
3. Generate a new plausible but uncommon attribute for this object that's not in the sentence.
4. Add the new attribute next to the selected object to make a new sentence.

The new sentence must meet the following three requirements:

1. The new sentence must describe a clearly new and different scene.
2. The new sentence must be fluent and grammatically correct.
3. The new sentence must make logical sense.

Here are some examples:

Original sentence: A large white airplane and a person on a lot.
Objects: ["airplane", "person"]
Selected object: airplane
New attribute: blue
New sentence: A large white and blue airplane and a person on a lot.

Original sentence: three people riding horses on a beach
Objects: ['three people', 'horses', 'beach']
Selected object: three people
New attribute: elderly
New sentence: Three elderly people riding horses on a beach.

(b) ADD-ATT.

Figure 9: Example prompt templates (black) and outputs (green) from ChatGPT for ADD hard negatives.

564 **Generating SWAP hard negatives.** To generate SWAP hard negatives, which do not require any
 565 new concepts, we simply prompt ChatGPT once with a temperature of 0.0. Unlike REPLACE, SWAP
 566 hard negatives are only possible when there are at least two atomic concepts of the same category,
 567 *i.e.*, either object or attribute. Thus, our prompt first queries ChatGPT whether it is possible to swap
 568 two atomic concepts in the input sentence to generate a new description. Only if the answer is yes,
 569 will ChatGPT then proceed to identify two swappable concepts and compose the corresponding new
 570 sentence by swapping the two concepts. Figure 8 shows the example templates and outputs for SWAP
 571 hard negatives.

572 **Generating ADD hard negatives.** Similar to the REPLACE, we also employ a three-step prompting
 573 procedure to generate ADD hard negatives. The only difference in the procedure is that we prompt

574 ChatGPT to add the generated new concept to the original caption, instead of using it to replace an
575 old concept. Figure 9 shows the example templates and outputs for ADD hard negatives.

576 D.3 Adversarial refinement

577 We detail the adversarial refinement procedure below. Given a text model M , we denote its output
578 score for the positive and negative caption of i -th image as $M(p_i)$ and $M(n_i)$. If $M(p_i) > M(n_i)$,
579 then the model could identify the correct caption for the i -th image without referring to it. For a test
580 set to be unattackable given the text model M , the expectation of M 's identifying the correct caption
581 should be as close to random guess as possible; in particular, we hope that $E_i[M(p_i) > M(n_i)] = 0.5$.
582 To achieve this for both the grammar model M_1 and plausibility model M_2 , we first calculate the score
583 difference $g_i^{(1)} = M_1(p_i) - M_1(n_i)$ and $g_i^{(2)} = M_2(p_i) - M_2(n_i)$, where the range of both $g^{(1)}$ and
584 $g^{(2)}$ is $[-1, 1]$. Then we split the 2D space of the joint range of $g^{(1)}$ and $g^{(2)}$ into 100×100 equal grids,
585 and for each pair of symmetric grids, e.g., $\{(g^{(1)}, g^{(2)}) | g^{(1)} \in (0.02, 0.04], g^{(2)} \in (-0.04, 0.06]\}$
586 and $\{(g^{(1)}, g^{(2)}) | g^{(1)} \in (-0.02, -0.04], g^{(2)} \in (0.04, -0.06]\}$, we preserve the same number of
587 data for both grids, therefore we ensure that for the resultant set, $E_i[M_1(p_i) > M_1(n_i)] = 0.5$ and
588 $E_i[M_2(p_i) > M_2(n_i)] = 0.5$.

589 D.4 Dataset information

590 We host SUGARCREPE on Github [13]. The data card [29] for SUGARCREPE, containing detailed
591 dataset documentation, is available at the dataset repository [14]. We provide a summary below.

592 **Dataset documentation.** SUGARCREPE is a benchmark for faithful vision-language compositionality
593 evaluation. Given an image, a model is required to select the positive text that correctly describes the
594 image, against another hard negative text distractor that differs from the positive text only by small
595 compositional changes. Each example consists of three fields:

- 596 • filename: The id to an image
- 597 • caption: Positive text correctly describing the image
- 598 • negative_caption: Hard negative text incorrectly describing the image

599 **Maintenance plan.** We are committed to maintain the dataset to address any technical issues. We
600 actively monitor issues in the repository.

601 **Licensing.** We license our work using MIT License [15]. All the source data we use is publicly released
602 by prior work [23].

603 **Author statement.** We the authors will bear all responsibility in case of violation of rights.

604 E Detailed evaluation results

605 E.1 Full evaluation results on existing benchmarks

606 We provide the full evaluation results over 17 pretrained CLIP models as well as 2 text-only models,
607 Vera [24] and the Grammar model [27], on existing compositionality benchmarks in Table 8. We see
608 that the text-only models, arguably without any vision-language compositionality, outperform most of
609 the pretrained CLIP models, achieving state-of-the-art performances on many benchmark tasks. This
610 implies that current benchmarks fail to faithfully reflect a model's vision-language compositionality.

¹³<https://github.com/RAIVNLab/sugar-crepe>

¹⁴https://github.com/RAIVNLab/sugar-crepe/blob/main/data_card.pdf

¹⁵<https://github.com/RAIVNLab/sugar-crepe/blob/main/LICENSE>

Table 8: Blind models (*i.e.*, Vera and Grammar model) outperform all 17 existing pretrained CLIP models on nearly all existing benchmark tasks. This implies that current benchmarks fail to faithfully measure a model’s vision-language compositionality.

Source	Model	CREPE			ARO				VL-Checklist		
		Atomic	Swap	Negate	VG-Relation	VG-Attribution	COCO-Order	Flickr30K-Order	Object	Attribute	Relation
Text-only model	Vera [24]	43.70	70.80	66.15	61.71	82.59	59.81	63.52	82.48	73.99	85.72
	Grammar [27]	18.15	50.88	9.77	59.55	58.38	74.33	76.26	57.95	52.35	68.50
OpenAI [30]	RN50	26.47	28.32	31.25	53.87	63.37	44.89	52.46	86.85	68.30	75.95
	RN101	27.63	32.74	12.50	52.43	62.93	29.86	39.34	86.44	67.93	71.75
	RN50x4	26.24	28.32	9.51	51.59	62.27	29.39	34.56	87.23	68.74	73.81
	ViT-B-32	22.31	26.55	28.78	51.12	61.33	37.14	47.18	87.00	68.80	77.04
	RN50x16	26.36	29.65	9.38	52.13	62.71	29.95	34.26	86.95	69.34	76.83
	RN50x64	26.82	30.09	23.57	51.00	62.56	40.54	46.74	87.71	68.61	74.97
	ViT-L-14	26.36	25.66	24.74	53.34	61.50	36.11	45.08	87.86	68.27	75.89
LAION [36]	ViT-H-14	23.70	25.22	16.54	50.33	62.93	25.79	30.96	85.39	68.46	71.13
	ViT-g-14	23.70	24.78	20.70	51.60	61.20	25.59	30.10	86.07	69.43	71.03
	ViT-bigG-14	23.58	24.78	17.97	51.61	61.89	25.24	30.22	84.66	67.80	66.48
	roberta-ViT-B-32	22.66	21.24	20.31	47.46	62.00	24.77	30.76	85.71	68.82	65.90
	xlm-roberta-base-ViT-B-32	21.16	20.80	12.76	47.93	59.73	23.85	30.32	86.06	70.41	63.01
	xlm-roberta-large-ViT-H-14	24.16	23.89	20.05	46.14	57.84	26.05	31.00	87.89	70.25	63.89
DataComp [12]	small: ViT-B-32	13.64	27.88	14.84	50.83	50.17	13.35	14.02	68.72	58.80	57.00
	medium: ViT-B-32	16.42	20.35	11.33	50.45	54.04	16.44	16.26	78.43	63.53	62.94
	large: ViT-B-16	18.15	17.26	17.06	48.82	53.21	21.49	26.44	84.73	65.72	64.81
	x-large: ViT-L-14	21.62	22.57	16.28	48.54	60.03	23.19	29.52	86.66	67.01	67.93

611 E.2 SUGARCREPE human evaluation

612 To compare the quality of the hard negatives generated in SUGARCREPE to those in current bench-
613 marks (*i.e.*, ARO+CREPE), we randomly sample 100 examples for each of the hard negative types:
614 REPLACE, SWAP, and NEGATE / ADD. Each example is organized to consist of (1) the original posi-
615 tive text, (2) its hard negative in ARO+CREPE, and (3) its hard negative in SUGARCREPE. For each
616 example, a human user rates whether the hard negative in ARO+CREPE or that in SUGARCREPE
617 is better (or tie) in terms of commonsense and grammatical correctness, respectively. Note that we
618 compare NEGATE in ARO+CREPE to ADD in SUGARCREPE, as both hard negatives are intended to
619 probe a model’s understanding of the *existence or not* of an atomic concept. Table 9 shows that hard
620 negatives in SUGARCREPE are much more sensical and fluent than that in ARO+CREPE across all
621 three different types. For instance, SUGARCREPE has 68% more sensical and 46% more fluent hard
622 negatives than ARO+CREPE on SWAP.

Table 9: Human evaluation results on the comparisons between hard negatives in ARO+CREPE and SUGARCREPE. We report the counts (out of 100 sampled examples) that the human user considers better or tie, w.r.t. both commonsense and grammatical correctness.

Hard-negative Type	Evaluation	Human counts of better examples		
		ARO+CREPE	SUGARCREPE	Tie
REPLACE	Commonsense	11	29	60
	Grammar	4	33	63
SWAP	Commonsense	4	68	28
	Grammar	4	46	50
NEGATE / ADD	Commonsense	1	26	73
	Grammar	1	35	64

274 References

- 275 [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson,
276 Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual
277 language model for few-shot learning. *Advances in Neural Information Processing Systems*,
278 35:23716–23736, 2022.
- 279 [2] Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora
280 Nozza, Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. Easily acces-
281 sible text-to-image generation amplifies demographic stereotypes at large scale. *arXiv preprint*
282 *arXiv:2211.03759*, 2022.
- 283 [3] Léon Bottou. From machine learning to machine reasoning. *Machine learning*, 94(2):133–149,
284 2014.
- 285 [4] Paola Cascante-Bonilla, Khaled Shehada, James Seale Smith, Sivan Doveh, Donghyun Kim,
286 Rameswar Panda, Gül Varol, Aude Oliva, Vicente Ordonez, Rogerio Feris, and Leonid Karlinsky.
287 Going beyond nouns with vision & language models using synthetic data, 2023.
- 288 [5] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz,
289 Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled
290 multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022.
- 291 [6] Jaemin Cho, Abhay Zala, and Mohit Bansal. Dall-eval: Probing the reasoning skills and social
292 biases of text-to-image generative transformers. *arXiv preprint arXiv:2202.04053*, 2022.
- 293 [7] Noam Chomsky and Morris Halle. Some controversial questions in phonological theory. *Journal*
294 *of linguistics*, 1(2):97–138, 1965.
- 295 [8] MJ Cresswell. *Logics and languages*. 1973.
- 296 [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, K. Li, and Li Fei-Fei. Imagenet: A large-
297 scale hierarchical image database. *2009 IEEE Conference on Computer Vision and Pattern*
298 *Recognition*, pages 248–255, 2009.
- 299 [10] Anuj Diwan, Layne Berry, Eunsol Choi, David Harwath, and Kyle Mahowald. Why is
300 winoground hard? investigating failures in visuolinguistic compositionality. *arXiv preprint*
301 *arXiv:2211.00768*, 2022.
- 302 [11] Sivan Doveh, Assaf Arbelle, Sivan Harary, Eli Schwartz, Roei Herzig, Raja Giryes, Rogerio
303 Feris, Rameswar Panda, Shimon Ullman, and Leonid Karlinsky. Teaching structured vision &
304 language concepts to vision & language models. In *Proceedings of the IEEE/CVF Conference*
305 *on Computer Vision and Pattern Recognition*, pages 2657–2668, 2023.
- 306 [12] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao
307 Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. Datacomp: In
308 search of the next generation of multimodal datasets. *arXiv preprint arXiv:2304.14108*, 2023.
- 309 [13] Madeleine Grunde-McLaughlin, Ranjay Krishna, and Maneesh Agrawala. Agqa: A benchmark
310 for compositional spatio-temporal reasoning. In *Proceedings of the IEEE/CVF Conference on*
311 *Computer Vision and Pattern Recognition*, 2021.
- 312 [14] Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and
313 Noah A. Smith. Annotation artifacts in natural language inference data. In *Proceedings of*
314 *the 2018 Conference of the North American Chapter of the Association for Computational*
315 *Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112, New
316 Orleans, Louisiana, June 2018. Association for Computational Linguistics.

- 317 [15] Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. Compositionality decomposed:
318 How do neural networks generalise? *Journal of Artificial Intelligence Research*, 67:757–795,
319 2020.
- 320 [16] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan
321 Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi,
322 Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021.
- 323 [17] Theo MV Janssen and Barbara H Partee. Compositionality. In *Handbook of logic and language*,
324 pages 417–473. Elsevier, 1997.
- 325 [18] Jingwei Ji, Ranjay Krishna, Li Fei-Fei, and Juan Carlos Niebles. Action genome: Actions as
326 compositions of spatio-temporal scene graphs. In *Proceedings of the IEEE/CVF Conference on*
327 *Computer Vision and Pattern Recognition*, pages 10236–10247, 2020.
- 328 [19] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie
329 Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting
330 language and vision using crowdsourced dense image annotations. *International journal of*
331 *computer vision*, 123(1):32–73, 2017.
- 332 [20] Ronan Le Bras, Swabha Swayamdipta, Chandra Bhagavatula, Rowan Zellers, Matthew Peters,
333 Ashish Sabharwal, and Yejin Choi. Adversarial filters of dataset biases. In *International*
334 *Conference on Machine Learning*, pages 1078–1088. PMLR, 2020.
- 335 [21] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: bootstrapping language-image
336 pre-training with frozen image encoders and large language models. In *ICML*, 2023.
- 337 [22] Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. BLIP: bootstrapping language-
338 image pre-training for unified vision-language understanding and generation. In Kamalika
339 Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors,
340 *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore,*
341 *Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 12888–
342 12900. PMLR, 2022.
- 343 [23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
344 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer*
345 *Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014,*
346 *Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- 347 [24] Jiacheng Liu, Wenya Wang, Dianshuo Wang, Noah A Smith, Yejin Choi, and Hannaneh
348 Hajishirzi. Vera: A general-purpose plausibility estimation model for commonsense statements.
349 *arXiv preprint arXiv:2305.03695*, 2023.
- 350 [25] Cewu Lu, Ranjay Krishna, Michael Bernstein, and Li Fei-Fei. Visual relationship detection
351 with language priors. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam,*
352 *The Netherlands, October 11–14, 2016, Proceedings, Part I 14*, pages 852–869. Springer, 2016.
- 353 [26] Zixian Ma, Jerry Hong, Mustafa Omer Gul, Mona Gandhi, Irena Gao, and Ranjay Krishna.
354 Crepe: Can vision-language foundation models reason compositionally? *arXiv preprint*
355 *arXiv:2212.07796*, 2022.
- 356 [27] John Morris, Eli Liland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. Textattack:
357 A framework for adversarial attacks, data augmentation, and adversarial training in nlp. In
358 *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing:*
359 *System Demonstrations*, pages 119–126, 2020.
- 360 [28] OpenAI. Chatgpt. 2022.

- 361 [29] Mahima Pushkarna, Andrew Zaldivar, and Oddur Kjartansson. Data cards: Purposeful and
362 transparent dataset documentation for responsible ai. In *2022 ACM Conference on Fairness,
363 Accountability, and Transparency*, pages 1776–1826, 2022.
- 364 [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-
365 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya
366 Sutskever. Learning transferable visual models from natural language supervision. In Marina
367 Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine
368 Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine
369 Learning Research*, pages 8748–8763. PMLR, 2021.
- 370 [31] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical
371 text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 372 [32] Arijit Ray, Filip Radenovic, Abhimanyu Dubey, Bryan A. Plummer, Ranjay Krishna, and
373 Kate Saenko. Cola: How to adapt vision-language models to compose objects localized with
374 attributes?, 2023.
- 375 [33] Yuval Reif and Roy Schwartz. Fighting bias with bias: Promoting model robustness by
376 amplifying dataset biases. *arXiv preprint arXiv:2305.18917*, 2023.
- 377 [34] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
378 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF
379 Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.
- 380 [35] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An
381 adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106,
382 2021.
- 383 [36] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman,
384 Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-
385 5b: An open large-scale dataset for training next generation image-text models. In *Thirty-sixth
386 Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- 387 [37] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba,
388 Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment
389 model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recogni-
390 tion*, pages 15638–15650, 2022.
- 391 [38] Harman Singh, Pengchuan Zhang, Qifan Wang, Mengjiao Wang, Wenhan Xiong, Jingfei Du,
392 and Yu Chen. Coarse-to-fine contrastive learning in image-text-graph space for improved
393 vision-language compositionality. 2023.
- 394 [39] Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela,
395 and Candace Ross. Winoground: Probing vision and language models for visio-linguistic
396 compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
397 Recognition*, pages 5238–5248, 2022.
- 398 [40] Junke Wang, Dongdong Chen, Zuxuan Wu, Chong Luo, Luowei Zhou, Yucheng Zhao, Yujia
399 Xie, Ce Liu, Yu-Gang Jiang, and Lu Yuan. Omnivl: One foundation model for image-language
400 and video-language tasks. *arXiv preprint arXiv:2209.07526*, 2022.
- 401 [41] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal,
402 Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. Image as a foreign language:
403 Beit pretraining for all vision and vision-language tasks. *arXiv preprint arXiv:2208.10442*,
404 2022.

- 405 [42] Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When
 406 and why vision-language models behave like bags-of-words, and what to do about it? In
 407 *International Conference on Learning Representations*, 2023.
- 408 [43] Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. Swag: A large-scale adversarial
 409 dataset for grounded commonsense inference. *arXiv preprint arXiv:1808.05326*, 2018.
- 410 [44] Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander
 411 Kolesnikov, and Lucas Beyer. Lit: Zero-shot transfer with locked-image text tuning. In
 412 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
 413 18123–18133, 2022.
- 414 [45] Tiancheng Zhao, Tianqi Zhang, Mingwei Zhu, Haozhan Shen, Kyusong Lee, Xiaopeng Lu,
 415 and Jianwei Yin. VI-checklist: Evaluating pre-trained vision-language models with objects,
 416 attributes and relations. *arXiv preprint arXiv:2207.00221*, 2022.

417 Checklist

- 418 1. For all authors...
- 419 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 420 contributions and scope? [Yes]
- 421 (b) Did you describe the limitations of your work? [Yes]
- 422 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 423 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 424 them? [Yes]
- 425 2. If you are including theoretical results...
- 426 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 427 (b) Did you include complete proofs of all theoretical results? [N/A]
- 428 3. If you ran experiments (e.g. for benchmarks)...
- 429 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 430 mental results (either in the supplemental material or as a URL)? [Yes]
- 431 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 432 were chosen)? [Yes]
- 433 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 434 ments multiple times)? [Yes]
- 435 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 436 of GPUs, internal cluster, or cloud provider)? [Yes]
- 437 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 438 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 439 (b) Did you mention the license of the assets? [Yes]
- 440 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 441 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 442 using/curating? [Yes]
- 443 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 444 information or offensive content? [N/A]
- 445 5. If you used crowdsourcing or conducted research with human subjects...
- 446 (a) Did you include the full text of instructions given to participants and screenshots, if
 447 applicable? [N/A]
- 448 (b) Did you describe any potential participant risks, with links to Institutional Review
 449 Board (IRB) approvals, if applicable? [N/A]
- 450 (c) Did you include the estimated hourly wage paid to participants and the total amount
 451 spent on participant compensation? [N/A]