

518 **A Appendix**

519 **A.1 Dataset Supplementary Materials**

- 520 1. Dataset documentation, metadata, and download instructions: anonymous.
521 2. Intended uses: we hope SeeTRUE will be used by researchers to evaluate image-text
522 matching models.
523 3. Author statement: We bear all responsibility in case of violation of right in using our
524 benchmark.
525 4. Each dataset’s license is described below. Our additional human annotations and gen-
526 erated images are licensed under CC-BY 4.0 license [https://creativecommons.org/
527 licenses/by/4.0/legalcode](https://creativecommons.org/licenses/by/4.0/legalcode).
528 5. Hosting & preservation: the dataset will be hosted in Huggingface Datasets, accessible and
529 available for open research.
530 6. SeeTRUE fields are presented in table 5.

531 Additional licensing details: we do publish the datasets we annotated in this work, and do not
532 re-publish the existing SNLI-VE and Winoground datasets. Full licenses:

- 533 1. MS COCO [28]: <https://cocodataset.org/#termsofuse>
534 2. EditBench [25]: <https://research.google/resources/datasets/editbench/>,
535 <https://www.apache.org/licenses/LICENSE-2.0>
536 3. DrawBench [5]: <https://imagen.research.google/>, [https://docs.google.com/
537 spreadsheets/d/1y7nAbmR4FREi6npB1u-Bo3GFdwdOPYJc617rB0xIRHY/edit#
538 gid=0](https://docs.google.com/spreadsheets/d/1y7nAbmR4FREi6npB1u-Bo3GFdwdOPYJc617rB0xIRHY/edit#gid=0)
539 4. Pick-a-Pick [29]: [https://huggingface.co/datasets/yuvalkirstain/pickapic_
540 v1](https://huggingface.co/datasets/yuvalkirstain/pickapic_v1)
541 5. SNLI-VE [27]: <https://github.com/necla-ml/SNLI-VE>
542 6. Winoground [24]: <https://huggingface.co/datasets/facebook/winoground>

Table 5: SeeTRUE Rows Examples

image	text	label	original_dataset_id	dataset_source
img1	A zebra to the right of a fire hydrant.	0	text_133_image_1228	drawbench
img2	A group of people standing next to bags of luggage.	1	text_105_image_1377	coco_t2i
img3	a tiny figurine is surrounded by cell phones on a table.	1	3786	editbench

543 **A.2 Human Annotation Process**



Figure 7: Annotation interface for determining whether a given image-text pair are aligned.

544 In order to provide reliable human labels for our datasets, we conducted an annotation process
545 using the SeeTRUE platform. The process comprised several steps, including setting qualification
546 requirements, providing instructions, and evaluating annotator agreement.

547 We set the basic requirements for our annotation task as follows: a percentage of approved assignments
 548 above 98%, more than 5,000 approved HITs, and annotator locations limited to the US, UK, Australia,
 549 or New Zealand. We selected 5 examples from our dataset for a qualification test and screened the
 550 annotators’ results. fig. 7 displays a sample of the Mechanical Turk user interface. The payment for
 551 the crowd-workers was 15-18 USD hourly.

552 The instructions provided were as follows:

553 *Evaluate the given image and text to determine if they match, selecting either “Yes” or “No”. Some*
 554 *images may be synthetically generated by a text-to-image model. To assess the match, mentally*
 555 *generate a textual description for the image (no need to write it down) and compare this generated*
 556 *description to the given text. If the descriptions closely resemble each other, mark “Yes”. If not, mark*
 557 *“No” and provide feedback on the specific issue causing the misalignment, focusing on the primary*
 558 *issue if multiple misalignments are present. If you encounter an image or text that may be offensive*
 559 *due to bias, race, NSFW content, etc., mark the checkbox to indicate this issue.*

560 Full agreement metrics are presented in table 6. As shown in the table, the percentage of cases
 561 where all annotators agreed and the Fleiss-Kappa scores vary across the datasets, with COCO-Con
 562 exhibiting the highest level of agreement and Drawbench the lowest. These differences highlight the
 563 varying levels of complexity within the datasets.

Table 6: Agreement metrics for different datasets.

Dataset	Full	Drawbench	COCO t2i	COCO-Con	PickaPic-Con
# Items	8,527	1,968	2,586	1,992	1,981
% all agreed	80	76	78	86	77
Fless-Kappa	0.72	0.66	0.68	0.81	0.69

564 A.3 Comparing VQ^2 variants

565 We have experimented with different variants of the VQ^2 pipeline.

Table 7: Comparing VQ^2 variants on all EditBench categories

Method	Extended answer candidates	Models used for validation	EditBench					
			Full	Color	Count	Material	Shape	Size
VQ^2 (NLI)	✗	VQA & NLI	66.8					
VQ^2 (NLI)	✓	VQA & NLI						
VQ^2 (PaLI comparison)	✓	VQA & VQA						
VQ^2 (yes/no)	✗	VQA						
VQ^2 (yes/no)	✓	VQA						

566 A.4 Reproducibility

567 To fine-tune BLIP2, we adjust only the Q-former parameters of the model using the Adam optimizer.
 568 We train the model for two epochs and designate 10% of the training set as a validation set for early
 569 stopping and use learning rate selection between $\{1e-5, 5e-5\}$. A single training took 5 hours on a
 570 linux server with one A6000 GPU. All experiments took <2 days.

571 Zero-shot VQ^2 : For 10,000 text-image pairs, the inference time of every step is as follows. Answer
 572 candidate generation: when using extended answer candidates – about 1 day. Otherwise, 12 hours.
 573 Question generation and filtering: When using extended answer candidates, about 2 days, otherwise,
 574 1 day. The last step only takes a few minutes.