# Contrasting intra-modal and ranking cross-modal hard negatives to enhance visio-linguistic compositional understanding

## Supplementary Material

## 6. Analysis of learned representations

In this section, we examine the image and text representations learned by our model. In particular, we investigate whether our method learns more distinct representations for positive and hard negative examples compared to those learned by CLIP and NegCLP. For each of CLIP, Neg-CLIP and CE-CLIP, we measure the intra-modal similarity between positive and hard-negative captions, as well as, the cross-modal similarity gap between positive and hard-negative image-caption pairs. We expect our method to reduce the intra-modal similarity and enlarge the cross-modal similarity gap compared to CLIP and NegCLIP. We report the results in Fig. 8, which shows that CE-CLIP achieves statistically significantly better intra-modal similarity (lower is better) and cross-modal similarity gap (higher is better) compared to CLIP and NegCLIP. To compute statistical significance, we used bootstrapping with 50,000 samples with confidence interval of 99%.
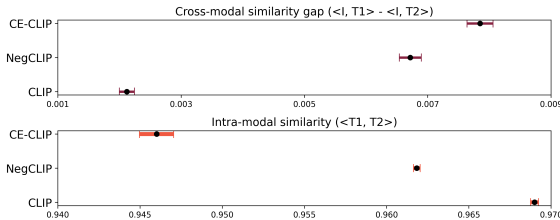


Figure 8. Analyzing the intra-modal similarities and cross-modal similarity gaps yielded by different methods on the ARO benchmark. T1 refers to positive caption while T2 refers to hard-negative caption. The red lines denote the standard errors obtained with bootstrapping 50,000 samples with confidence interval of 99%.

## 7. Qualitative Examples

Fig. 9 and Fig. 10 illustrate some side-by-side comparisons of hits and misses by CE-CLIP versus NegCLIP.

## 8. Benchmark

The statistics of benchmarks we use are shown in Tab. 6.

| Benchmark | Task | # image-text pairs |
|---|---|---|
| ARO-Relation | Relation | 24k |
| ARO-Attribution | Attribution | 28.7k |
| VALSE | Linguistic Phenomena | 6.8k |
| VL-CheckList | Objects, Attributes and Relations | 410k |
| SugarCrepe | Objects, Attributes and Relations | 7.5k |

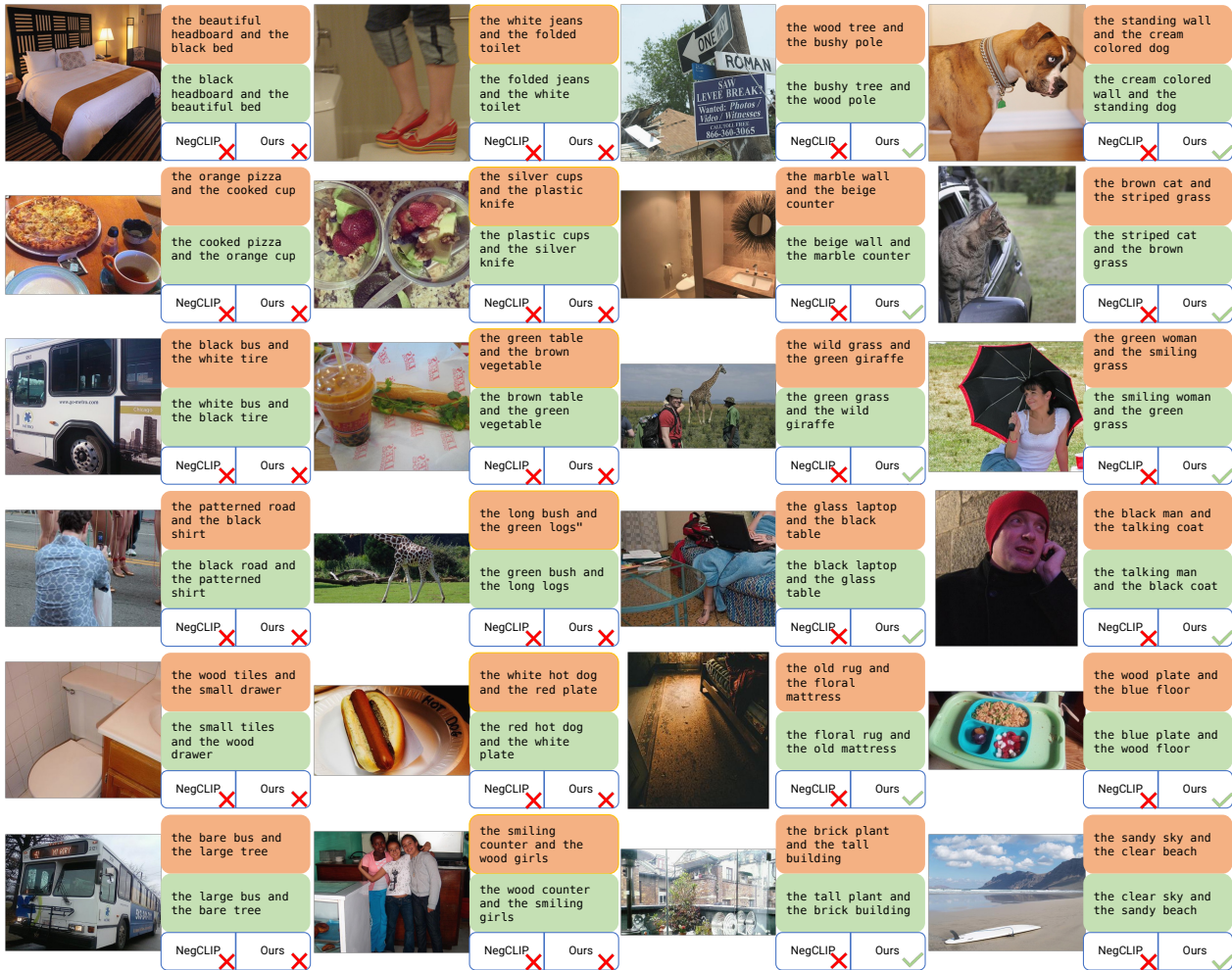Table 6. Overview of vl-compositional benchmarks.

## 9. Acknowledgement

Figure 9. Some qualitative examples from ARO-Attribute. Caption in red box is unmatched and in green box is matched. ✓ represents model predicates correctly and ✗ means wrong.
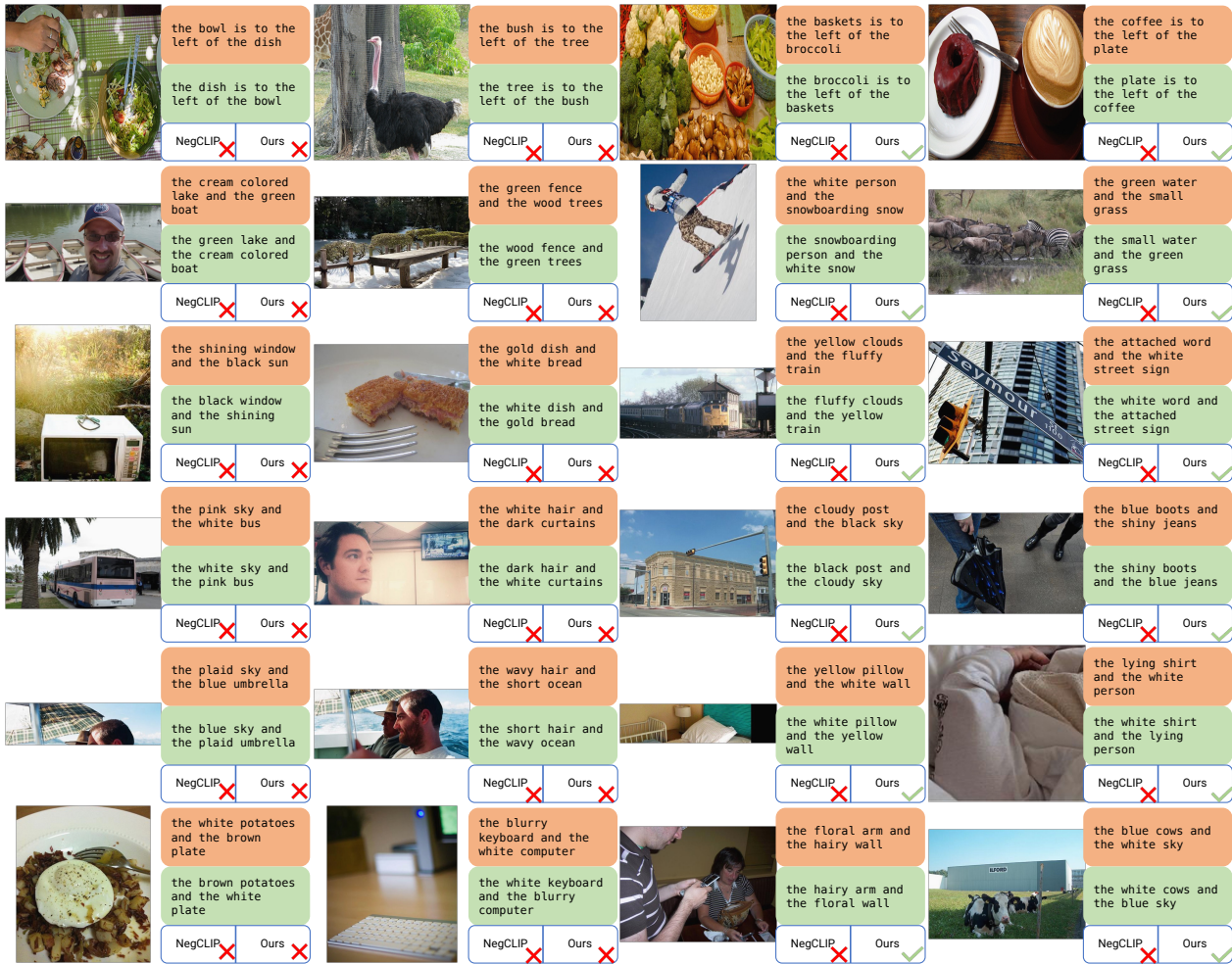
Figure 10. Some qualitative examples from ARO-Relation. Caption in red box is unmatched and in green box is matched. ✓ represents model predicates correctly and ✗ means wrong.