





Figure 2. Qualitative examples on coarse-grained and fine-grained counterfactual samples.

## 2.2. Qualitative Analysis

We provide more coarse-grained and fine-grained qualitative examples, visualized in Figure 2. These images are from MS-COCO and have been trained by our model, while the text queries are newly annotated. We show our model’s classification performance on coarse-grained and fine-grained counterfactual samples in two columns, respectively.

Specifically, the left samples are random image-text pairs without any semantic connections. The right samples are matched image-text pairs with the attribute words in texts changed into counterfactual words. Our model successfully identifies all the coarse-grained samples and most fine-grained samples. This result indicates that training a C-REC model on fine-grained samples also contributes to detecting coarse-grained samples, thus covering almost all counterfactual situations.

## References

- [1] Yutao Hu, Qixiong Wang, Wenqi Shao, Enze Xie, Zhenguo Li, Jungong Han, and Ping Luo. Beyond one-to-one: Rethinking the referring image segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4067–4077, 2023. 1
- [2] Chang Liu, Henghui Ding, and Xudong Jiang. Gres: Generalized referring expression segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23592–23601, 2023. 1