

supplementary

October 17, 2024

1 Examples of failed cases:

The example shows that while both models retrieved correct images—i.e., images of streets—one of the images retrieved by SR_{clip} was marked incorrect during evaluation because it lacked "street" in its labels. Since we rely on the provided labels for evaluating model performance, this led to a mismatch. There are several similar cases in the Conceptual Captions dataset, where images contain specific objects or features that are missing from the labels. This is likely due to the large number of labels in the dataset, inaccuracies in the word mapping, and the fact that the labels are not human-generated.

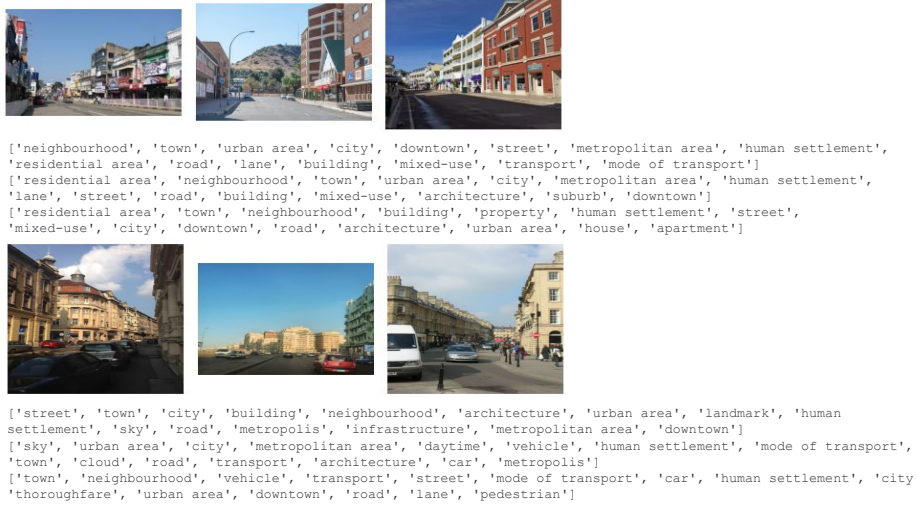


Figure 1: Retrieval images and their labels for Conceptual Captions dataset for Query "street without trees" using CLIP and SR_{clip}