
Supplementary: Pick-a-Pic: An Open Dataset of User Preferences for Text-to-Image Generation

Anonymous Author(s)
Affiliation
Address
email

1 Appendix

2 Training a Scoring Function

3 We further explored an alternative loss function that is closer to CLIP’s original objective. In this loss
4 function, we also incorporate the remaining examples in the batch as in-batch negatives. To elaborate,
5 considering the notation established in Section 3 and y_1^k, y_2^k denoting the images corresponding to
6 the k -th example in the batch, we formulate \hat{p}^k as follows:

$$\hat{p}_i^k = \frac{\exp s(x, y_i)}{\sum_k \sum_{j=1}^2 \exp s(x, y_j^k)} \quad (1)$$

7 We anticipated that this objective function would maintain the general capabilities of CLIP with
8 minimal loss in performance. However, our findings demonstrated that PickScore significantly
9 outperforms this objective function, as the latter only produced a scoring function that achieves an
10 accuracy of 65.2 on the Pick-a-Pic test set.