Tagging like Humans: Diverse and Distinct Image Annotation (Supplementary Material)

Baoyuan Wu¹, Weidong Chen¹, Peng Sun¹, Wei Liu¹, Bernard Ghanem², and Siwei Lyu³

¹Tencent AI Lab ²KAUST ³University at Albany, SUNY

1. Qualitative Results

Here we present some qualitative results of to compare the performance of DIA [3] and our proposed method D²IA-GAN. The results on ESP Game [2] data are shown in Fig. 3, and those on IAPRTC-12 [1] data are shown in Fig. 4 (see the last page). In each sub-figure, the left is the original image, while the right shows the ground-truth complete tag list and the annotation results of DIA and $D^{2}IA$ -GAN. For each method, we present two cases, including 3 tags (*i.e.*, each single tag subset includes at most 3 tags) and 5 tags (i.e., each single tag subset includes at most 5 tags). At each case, we present at most 3 tag subsets. As described in Section 5.5 of the main manuscript, the DPP sampling process in each method (DIA or D^2 IA-GAN) is run 10 times to obtain 10 tag subsets. Then 3 subsets with the largest tag weight summations are picked as the presented results. If there are same subsets in these 3 subsets, the redundant subsets are removed. Furthermore, we also present the F_{1-sn} score of the ensemble subset of these picked subsets. From these qualitative results, we can see that (1) in each single tag subset of both DIA and D²IA-GAN, the included tags are semantically distinct to each other; (2) the tag subsets of $D^{2}IA$ -GAN are more diverse than those of DIA, so the corresponding ensemble tag subset of D²IA-GAN covers more semantic meanings than that of DIA.

2. Analysis of Human Annotations

As demonstrated in Introduction of the main manuscript, we have conducted a human annotation experiment by asking 3 human annotators to independently annotate the first 1000 test images in the IAPRTC-12 dataset, with the requirement of 'describing the main contents of one image using as few tags as possible. Here we present some analysis about these human annotation results.

The first analysis focuses on the single tag subset from



Figure 1. Statistics of the sizes of single tag subsets of 3 human annotations, on the first 1000 test images of IAPRTC-12 [1].



Figure 2. Statistics of the sizes of ensemble tag subsets of human annotations, DIA [3] and D^2 IA-GAN, on the first 1000 test images of IAPRTC-12 [1].

each human annotator. As shown in Fig. 1, the statistics of the sizes of single tag subsets of 3 human annotators are plotted in 3 sub-figures separately. The average sizes of single tag subsets from 3 human annotations are 3.99, 5.15, 4.43, respectively. This demonstrates that the single human annotator tends to describe the image content using a few **relevant** and **distinct** tags.

The second analysis is about the ensemble tag subset, derived by merging 3 human annotated tag subsets for the

same image and removing the repeated tags. The statistics of the sizes of ensemble tag subsets of human annotations are presented as the green bar in Fig. 2. The average size of 1000 human annotated ensemble tag subsets is 11.22. The gap between the average size of single tag subsets and that of ensemble tag subsets reveals that different human annotators tend to give different tag subsets, *i.e.*, diverse. Moreover, we also present the statistics of the sizes of ensemble tag subsets produced by DIA [3] and our proposed method D²IA-GAN in Fig. 2, distinguished by the yellow and purple colors, respectively. The ensemble tag subset is derived by merging 5 tag subsets produced by DIA or D²IA-GAN and removing the repeated tags. These 5 subsets are exactly the results evaluated in our experiments (see Section 5.2 in the main manuscript), and each subset includes at most 5 tags. Specifically, the sizes of ensemble tag subsets of DIA range from 5 to 9, and the average size is 7.09. In contrast, the size range of D^2 IA-GAN is from 3 to 18, and the average size is 7.95. This comparison indicates that the diversity between the tag subsets produced by D²IA-GAN is larger than that by DIA, and the ensemble tag subset of $D^{2}IA$ -GAN can cover more semantic meanings than that of DIA. However, the gap between the average size of ensemble tag subsets produced by human annotations and that by $D^{2}IA$ -GAN reminds us that the diversity of our automatic tagging results is still smaller than the diversity of human annotations. The main reason is that both the training and inference of D²IA-GAN are constrained in a fixed set of candidate tags provided by the dataset (291 candidate tags in IAPRTC-12, and 268 candidate tags in ESP Game [2]). In contrast, the tags used by human annotators are unconstrained, and many are out of the range of the fixed set we used. It is expected that D²IA-GAN could produce more human-like tags if we train it on a larger set of candidate tags derived from the collection of real human annotated tags. This is our future research work.

References

- M. Grubinger, P. Clough, H. Müller, and T. Deselaers. The iapr tc-12 benchmark: A new evaluation resource for visual information systems. In *International Workshop OntoImage*, pages 13–23, 2006. 2
- [2] L. Von Ahn and L. Dabbish. Labeling images with a computer game. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 319–326. ACM, 2004. 2, 3
- [3] B. Wu, F. Jia, W. Liu, and B. Ghanem. Diverse image annotation. In CVPR, 2017. 2



Figure 3. Some qualitative results on ESP Game. The value in brackets at the end of each row indicates the F_{1-sp} score of the ensemble subset of the subsets in the same row.



Figure 4. Some qualitative results on IAPRTC-12. The value in brackets at the end of each row indicates the F_{1-sp} score of the ensemble subset of the subsets in the same row.