

Supplementary Materials for Understanding and Evaluating Racial Biases in Image Captioning

Dora Zhao Angelina Wang Olga Russakovsky
Princeton University

{dorothy, angelina.wang, olgarus}@cs.princeton.edu

A. Comparing collected gender annotations with automatically derived ones

We explore extending the schema introduced in Sec. 3.2 for deriving gender labels from captions in three ways: 1) only labeling images where there is a `person` who has a bounding box greater than 5,500 pixels, 2) expanding the list of gendered words beyond “man” and “woman”, and 3) having different cutoffs for how many captions (of the 5 per image) need to mention a gender for the image to be labeled. We call the use of the gendered set {man, woman} “few,” and that of our expanded set “many.”

Our expanded set “many” consists of the following words: [“male”, “boy”, “man”, “gentleman”, “boys”, “men”, “males”, “gentlemen”] and [“female”, “girl”, “woman”, “lady”, “girls”, “women”, “females”, “ladies”].

Our results in Fig. 1 show that while these extensions significantly increase both the number of images correctly labeled and the accuracy of labeled images, all methods are inaccurate and/or incomplete. The gender labels derived from captions remain highly imperfect, as expected, cautioning against automated means of gender derivation [40].

B. Racial descriptors

When searching for descriptors of race and ethnicity in Sec. 4.1, we first convert the captions to lowercase. We then use the following keywords [“white”, “Caucasian”, “Black”, “African”, “Asian”, “Latino”, “Latina”, “Latinx”, “Hispanic”, “Native”, and “Indigenous”] — also in lowercase — to query the captions.

C. Caption performance

In Sec. 4.2 we assess the differences in caption performance for BLEU [55], CIDER [68], and SPICE [2] when evaluated on the COCO 2014 validation set. We extend this analysis by providing the overall scores across the four image captioning models and looking at two additional automated image captioning metrics.

To start, we look at the performance for our four models.

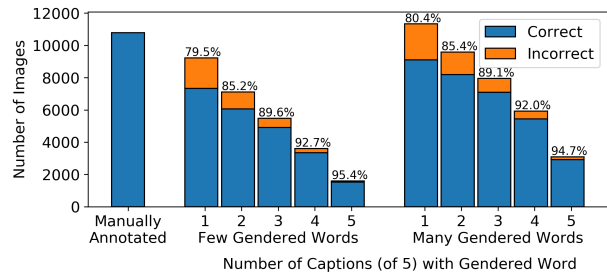


Figure 1: Comparison of various ways of automatically deriving image-level gender annotations from existing image captions. “Few” and “Many” gendered words refers to the size of the set of gendered words considered, and numbers refer to required captions that mention a gendered word, higher numbers limiting the images that can be labeled. Percentage over each bar indicates accuracy. All methods are imperfect and noisy, cautioning against the use of automatically-derived gender annotations.

As seen in Tbl. 1, **Oscar** has the best performance across all metrics. Further, we see that the newer transformer-based models outperform older models (e.g. **FC**, **Att2in**, and **Dis-cCap**) across all metrics as well.

We also report the differences in performance between lighter and darker images for two commonly used image captioning metrics — METEOR [5] and ROUGE [46] (Tbl. 2). Similar to the results for BLEU and CIDER, the differences for METEOR and ROUGE are greater for **Att2in**, **Transformer**, and **Oscar**. **AoANet** also shows some slight differences in performance for METEOR and ROUGE. This supports our observation that the better performing captioning models also tend to show greater discrepancies in performance between lighter and darker images.

Table 1: The caption performance as measured by BLEU [55], CIDEr [68], and SPICE [2] multiplied by 100 on the COCO 2014 validation dataset. Error bars represents 95% confidence intervals across random seeds used to train 5 models per architecture.

	BLEU	CIDEr	SPICE
FC [59]	28.2 ± 0.4	87.2 ± 0.4	16.9 ± 0.2
Att2in [59]	31.5 ± 0.2	94.2 ± 0.3	18.4 ± 0.1
DiscCap [50]	23.7 ± 0.7	71.1 ± 2.3	18.8 ± 0.2
Transformer [67]	33.9 ± 0.6	97.2 ± 2.2	20.2 ± 0.2
AoANet [33]	38.7 ± 1.6	116.2 ± 4.6	22.2 ± 0.6
Oscar [45]	40.0 ± 0.4	120.0 ± 0.5	23.0 ± 0.3

D. Vocabulary differences coefficients

In Sec. 4.4 we explore the different word choices in the captions describing `lighter` and `darker` images. We provide the most predictive words for `lighter` and `darker` across the manual captions and the automatically generated captions in Tbl. 3.

