Figure 1: Left: visualisation of feature channels. The number on the right top corner is the channel number. The word that has the highest correlation $\alpha_{i,j}$ in Eqn.1 with the channel is shown under the image. Right: ablation study.
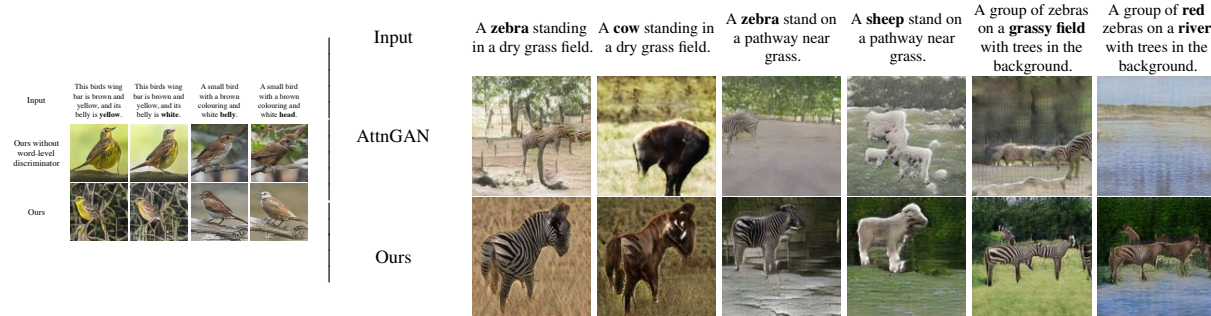


Figure 2: Left: ablation study. Right: Qualitative comparison of two methods on the COCO dataset. Odd-numbered columns show the original text and even-numbered ones the modified text. Please zoom in to see clearer.

1  **R1: The effectiveness of channel-wise attention.** We experimentally find that the channel-wise attention highly
2  correlates with semantically meaningful parts while the spatial attention focuses on colour descriptions. Fig.1(L) shows
3  several channels of feature maps transformed with our channel-wise attention. It shows that channels closely correlate
4  to semantic parts, such as belly, eye, wing, crown. This phenomenon is further verified by the results shown in Fig.1(R).
5  Without channel-wise attention, the model fails to generate controllable results when we modify the text related to parts
6  of a bird. In contrast, our full model with channel-wise attention can generate much better controllable results.

7  **R1: The necessity of attention.** In the paper Fig. 6 and Tab. 2, *SPM only* means the model does not incorporate the
8  channel-wise attention but still has the spatial attention. We will clarify this in our main paper. Fig. 6 *SPM only* verifies
9  that spatial attention works well if only colour information is modified. Our channel-wise attention is responsible for
10  parts related modification as described above.

11  **R1: Word-level discriminator.** The word-level discriminator can help better disentangle different visual attributes and
12  facilitate image manipulation as shown in Fig.2(L). With word-level discriminator, our model achieves better results,
13  e.g., the original shape and colour of the bird are well preserved.

14  **R1: Novelty of SPM.** To our best knowledge, we are the first to apply perceptual loss in controllable text-to-image
15  generation and show its effectiveness on reducing randomness involved in generation, which makes manipulation stable.

16  **R2: Channel-wise attention.** Please refer to "R1: The effectiveness of channel-wise attention."

17  **R2: The importance of self-attention.** We adopt the pre-trained RNN used in AttnGAN. The objective function
18  used to train RNN is to improve text-image matching score based on cosine similarity. Thus, it is reasonable to use
19  self-attention to reduce impact from less important words which has small cosine similarity.

20  **R2: SPM and diversity of the model.** The SPM would not affect the diversity of the model, as the diversity comes
21  from the random noise sampled from a normal distribution in stage 0. We will clarify this in our paper.

22  **R2: The controllability and out-of-distribution on COCO.** In this rebuttal, we focus on the animal subset of COCO.
23  As shown in Fig.2, our model can effectively manipulate the images compared with AttnGAN and also work well on
24  out of distribution queries, e.g., red zebras on the river.

25  **R3: Shape of the visual features, sensitivity.** $v$ in equation 1 indicates the generated image features whose shape is
26  $B \times N \times (H * W)$, where $N$ is 32, $H * W$ is $64 * 64$ in stage 1 and $128 * 128$ in stage 2. The $v$ in equation 2 indicates
27  real image features whose shape is same as generated images. We will clarify this.

28  **R3: CUB and COCO experiments** COCO has much more diversified descriptions where each only contains few
29  examples. The low density of corresponding images is a reason for the different behaviours on the two datasets. Also,
30  captions in COCO are more abstract and focus on the category of objects, which makes text-to-image generation be
31  more challenging on COCO. However, we still produce much better results compared with AttnGAN shown in Fig.2(R).