

A text to image generation (T2I) model aims to generate photo-realistic images which are semantically consistent with the text descriptions. Built upon the recent advances in generative adversarial networks (GANs), existing T2I models have made great progress. However, a close inspection of their generated images reveals two major limitations: (1) The condition batch normalization methods are applied on the whole image feature maps equally, ignoring the local semantics; (2) The text encoder is fixed during training, which should be trained with the image generator