

Project D - TokenCompose: Text-to-Image Diffusion with Token-level Supervision

Julien DELAVANDE, Soël MEGDOUD

08/01/2025

1 Introduction

Recent advancements in text-to-image generation models, such as Latent Diffusion Models (LDMs), have achieved impressive results in generating high-quality, photorealistic, and diverse images conditioned on textual prompts. However, a significant challenge remains: ensuring consistent alignment between the content of user-specified text prompts and the images generated by these models. This problem becomes particularly pronounced when generating images involving multiple object categories, especially configurations uncommon in real-world scenarios. Standard LDMs often fail to adequately compose multiple objects, leading to missing or poorly arranged instances in the generated images.

TokenCompose[3] addresses this limitation by introducing token-wise consistency terms during the finetuning stage of Stable Diffusion. These terms leverage pretrained image understanding models, such as Segment Anything (SAM), to provide token-level supervision. The proposed method enhances the model’s ability to accurately compose multiple object categories in generated images, significantly improving both object accuracy and photorealism.

2 Latent Diffusion Models & proposed improvements

Latent Diffusion Models operate by encoding an image into a latent space using a variational autoencoder (VAE). During training, random noise is added to the latent representation, and the denoising function, parameterized by a U-Net architecture, is optimized to predict and remove the noise. To condition image generation on textual prompts, text embeddings are injected into the U-Net layers via cross-attention. The training process involves computing the loss between the predicted noise and the ground truth noise.

$$\mathcal{L}_{\text{LDM}} := \mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(y))\|_2^2 \right] \quad (1)$$

Where ϵ is the random noise, ϵ_{θ} is a neural network, $z_t = E(x_t)$ is the encoded image by the VAE and $\tau_{\theta}(y)$ is the text encoded through CLIP[2].

Despite these advances, existing LDMs optimize only for denoising, which fails to explicitly align text tokens with image regions. This limitation results in poor token-level understanding and suboptimal composition of multiple objects during inference.

2.1 Token-Level Attention Loss ($\mathcal{L}_{\text{token}}$)

To improve token-level grounding, TokenCompose introduces a token-level attention loss, $\mathcal{L}_{\text{token}}$, which supervises the cross-attention activations between text tokens and corresponding regions in the image. The method uses binary segmentation maps M_i for each text token i , generated automatically by pretrained segmentation models.

The token-level attention loss encourages activations in the cross-attention maps to focus on the target regions defined by B_i , normalized by the total activations. This approach ensures that each token’s embedding aligns with its corresponding region in the image. The loss function for $\mathcal{L}_{\text{token}}$ is defined as:

$$\mathcal{L}_{\text{token}} = \frac{1}{N} \sum_{i=1}^N \left(1 - \frac{\sum_{u \in B_i}^{L_{z_t}} A_{(i,u)}}{\sum_u^{L_{z_t}} A_{(i,u)}} \right)^2 \quad (2)$$

where $A_{(i,u)}$ represents the scalar attention activation at a spatial location u and B_i are predicted spatial regions $B_i = \{u \in M_i | u = 1\}$.

2.2 Pixel-Level Attention Loss ($\mathcal{L}_{\text{pixel}}$)

While $\mathcal{L}_{\text{token}}$ improves token-level alignment, it may lead to overly concentrated activations within certain subregions of the target areas. To address this issue, TokenCompose introduces a pixel-level attention loss, $\mathcal{L}_{\text{pixel}}$, which applies a binary cross-entropy objective to constrain activations at the pixel level. This ensures that activations are distributed appropriately across the target regions defined by the segmentation maps.

$$\mathcal{L}_{\text{pixel}} = -\frac{1}{L_{\tau_{\theta}(y)} L_{z_t}} \sum_{i=1}^{L_{\tau_{\theta}(y)}} \sum_{u=1}^{L_{z_t}} \left[\mathcal{M}_{(i,u)} \log(A_{(i,u)}) + (1 - \mathcal{M}_{(i,u)}) \log(1 - A_{(i,u)}) \right] \quad (3)$$

2.3 Loss for training

The combined optimization objective during training is:

$$\mathcal{L}_{\text{TC}} = \mathcal{L}_{\text{LDM}} + \sum_{m=1}^M \left(\lambda \mathcal{L}_{\text{token grounding}}^{(m)} + \gamma \mathcal{L}_{\text{pixel}}^{(m)} \right) \quad (4)$$

where λ and γ are scaling factors to balance the contributions of the token- and pixel-level losses.

2.4 Effects of token level supervision

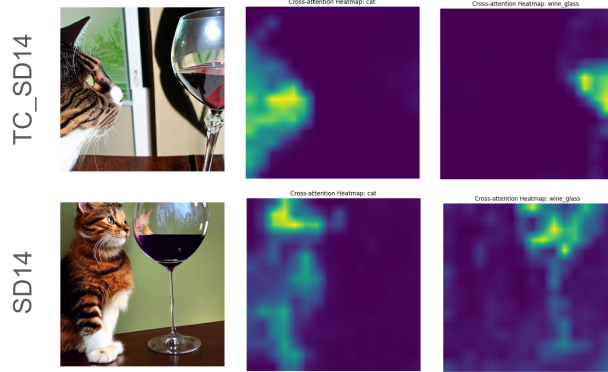


Figure 1: Generated images for "A cat and a wine glass" prompt by Stable Diffusion and TokenCompose

Here of attention maps for 'cat' and 'wine glass' tokens for an image generated by Stable Diffusion and TokenCompose in our tests. The attention is much more isolated on the objects for TokenCompose, especially for the glass. In addition, attention is more evenly distributed across the object thanks to Pixel level loss.

3 Multi-Category Instance Composition Evaluation Methodology

We evaluated the compositional capabilities of the models using the MULTIGEN benchmark, which examines the generation of multiple object categories within a single image based on textual prompts. MULTIGEN poses a challenging task by randomly sampling 5 distinct categories (e.g., "A, B, C, D, and E") and formatting them into a sentence (e.g., "A photo of A, B, C, D, and E"). These prompts are used as input conditions for text-to-image diffusion models to generate corresponding images.

To evaluate performance, we employ a robust open-vocabulary detector to identify the presence of the specified categories in the generated images. In the original paper, for each dataset (COCO and

ADE20K), 1,000 prompts are sampled to generate 10 images per prompt, resulting in a total of 10,000 generated images per dataset. Then, the success rate of generating 2 to 5 specified categories out of 5 (MG2-5) for each round as well as the standard deviation across the 10 rounds is calculated. In our implementation, we only randomly chose 50 prompts and generated 5 images per prompts to limit computational time.

4 Proposed CLIP Fine-Tuning Enhancement

In the original framework, the text encoder utilized in the model is derived from CLIP, with its weights frozen during training. We hypothesized that fine-tuning CLIP for the task could yield improved results, particularly for generating images containing multiple object categories in varied spatial arrangements.

CLIP was originally trained on high-quality images, where objects were often centrally positioned and clearly depicted. However, this central bias could limit its ability to represent objects in complex compositions. To address this, we leveraged the COCO dataset and its segmentation masks to create a fine-tuning dataset for CLIP. By extracting object instances from images, including cases where objects appear in non-central positions or are occluded by others, we introduced variability to the training data.

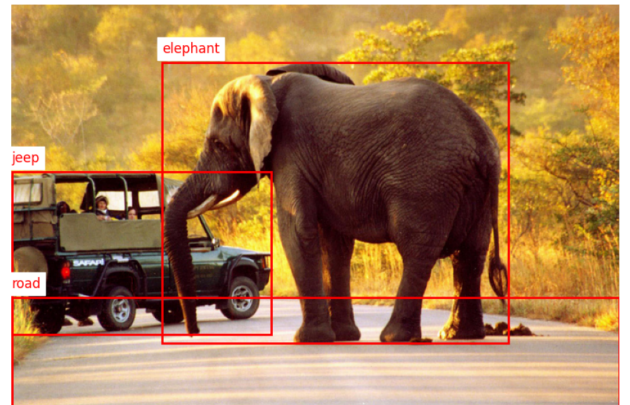


Figure 2: Bounding boxes extracted from segmentation maps in Coco dataset

The enhanced data set was then used to fine-tune CLIP using a contrastive learning objective[1].

Contrastive learning involves aligning semantically similar pairs (e.g., a textual description of an object and its corresponding image segment) while pushing apart embeddings of dissimilar pairs. The InfoNCE

loss, employed for this purpose, is defined as:

$$\mathcal{L}_{\text{InfoNCE}} = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\text{sim}(z_i, z_i^+)/\tau)}{\sum_{j=1}^N \exp(\text{sim}(z_i, z_j)/\tau)}, \quad (5)$$

where z_i and z_i^+ represent the embeddings of a positive pair (e.g., a segment and its corresponding textual token), $\text{sim}(\cdot, \cdot)$ is a similarity function (e.g., cosine similarity), and τ is a temperature parameter controlling the distribution sharpness.

This fine-tuning process aims to bring the embeddings of object representations closer to their corresponding extracted images, while ensuring that unrelated embeddings remain distant.

5 Results & Conclusion

5.1 Results

After finetuning, we evaluated the TokenCompose model based on Stable Diffusion 1.4 with the fine-tuned CLIP model as a text encoder. The MULTIGEN benchmark has been set with 50 prompts x 5 images for computational reasons.

Model	MG ₂	MG ₃	MG ₄	MG ₅
TC.SD14	93.50 _{0.77}	74.00 _{3.64}	31.50 _{6.88}	3.00 _{1.63}
SD14	84.50 _{2.18}	43.50 _{5.62}	9.50 _{5.62}	1.00 _{1.00}
TC.SD14 + CLIP	94.50 _{2.08}	71.00 _{3.22}	27.00 _{2.14}	2.00 _{1.31}
TC14 Original Paper	98.08 _{0.40}	76.16 _{1.04}	28.81 _{0.95}	3.28 _{0.48}
SD14 Original Paper	90.72 _{1.33}	50.74 _{9.59}	11.68 _{0.45}	0.88 _{0.21}

Table 1: MULTIGEN Benchmark Results: Success rates (%) and standard deviations for generating multiple categories. The results from the original paper are included for comparison.

As you can see, our fine tuning downgrade performances, especially from for 4 and 5 objects composition generation (MG_4 and MG_5) and performance remain quite similar for 2 and 3 objects composition. Unfortunately, the fine tuning is not an improvement.

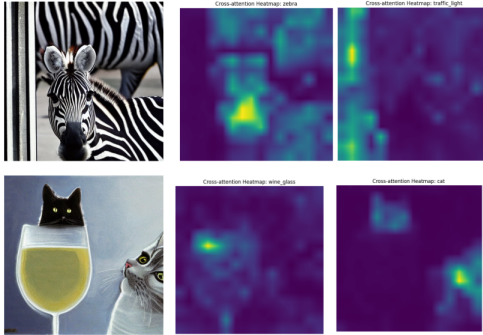


Figure 3: Generated images and attention maps. Prompts: "A zebra next to a traffic light" and "A cat and a wine glass".

Attention maps are also degraded, with attention less well distributed on the objects to be generated.

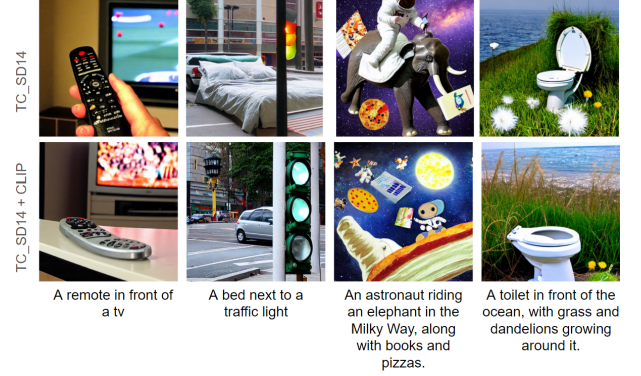


Figure 4: Examples of generated images with our fine-tuned model and TokenCompose

5.2 Conclusion

Our experiments confirm the original article: adding loss during training focus attention activations to specific areas of the image, making it easier to compose images with several objects, as demonstrated by the MULTIGEN benchmark.

Unfortunately, our proposition does not improve performance and several reasons could explain it. Some of the labels on the masks are not precise enough and some of the bounding boxes are sometimes badly positioned which can be detrimental to training. Unfortunately, we believe that this finetuning downgraded CLIP. In addition, we were unable to re-train the entire model (Stable Diffusion) due to hardware constraints and TokenCompose was trained on standard CLIP. One way of improving this is to re-train the model with our version of CLIP to see the difference.

References

- [1] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations, 2020.
- [2] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021.
- [3] Zirui Wang, Zhizhou Sha, Zheng Ding, Yilin Wang, and Zhuowen Tu. Tokencompose: Text-to-image diffusion with token-level supervision, 2024.