

Null-text Inversion for Editing Real Images using Guided Diffusion Models

Ron Mokady*^{1,2}, Amir Hertz*^{1,2}, Kfir Aberman¹, Yael Pritch¹, and Daniel Cohen-Or^{† 1,2}

¹Google Research

²The Blavatnik School of Computer Science, Tel Aviv University

MI2 Seminar

Dawid Płudowski

December 2nd, 2024



Authors



Ron Mokady*



Amir Hertz*



Kfir Aberman



Yael Pritch



Daniel Cohen-Or

Google research
600+ citations

Conference on Computer Vision and Pattern Recognition 2023





Motivation

To edit a real image using these state-of-the art tools, one must first invert the image with a meaningful text prompt into the pretrained model's domain. In this paper, we introduce an accurate inversion technique and thus facilitate an intuitive text-based modification of the image.





Background I

In 2022, main authors released “*Prompt-to-Prompt Image*

- ▶ *Editing with Cross Attention Control* which allows to make good edits of images

- ▶ A big issue of this work is the need to inverse the real image; in the article, they were focused on artificial images mostly



When we need to show the method works

Prompt-to-Prompt, 2022

“A black bear is walking in the grass.”



real image



reconstructed

“Landscape image of trees in a valley...”



real image



reconstructed

When we need to show the method must be improved

NULL-text optimization, 2022

Input Image



DDIM Inversion



Input caption: “Zoom photo of flower”



Input caption: “A cat sitting next to another cat”



Wydział Matematyki i Nauk Informacyjnych

POLITECHNIKA WARSZAWSKA

MI DATALAB



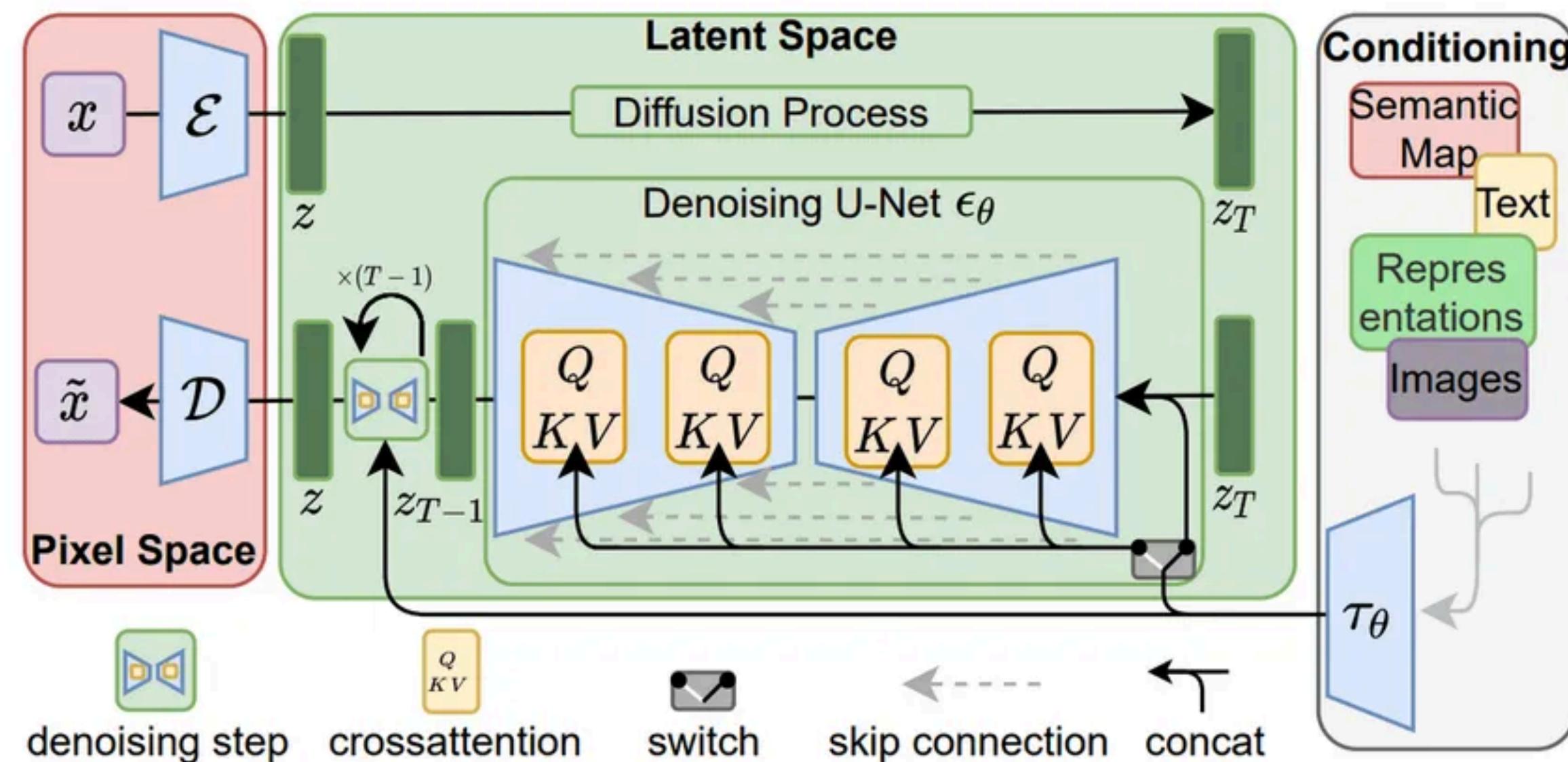
Background II

- Stable diffusion (autoencoder architecture)
- Classifier-free guidance



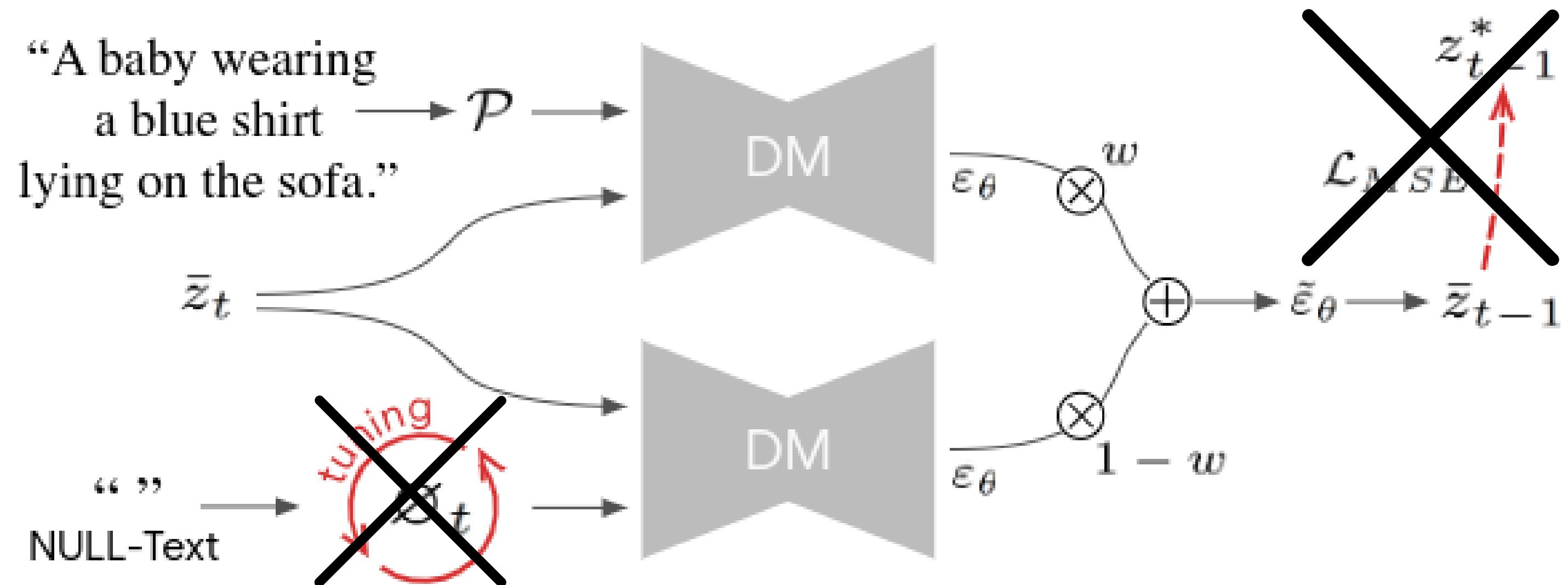


Reminder - stable diffusion





Reminder - classifier-free guidance



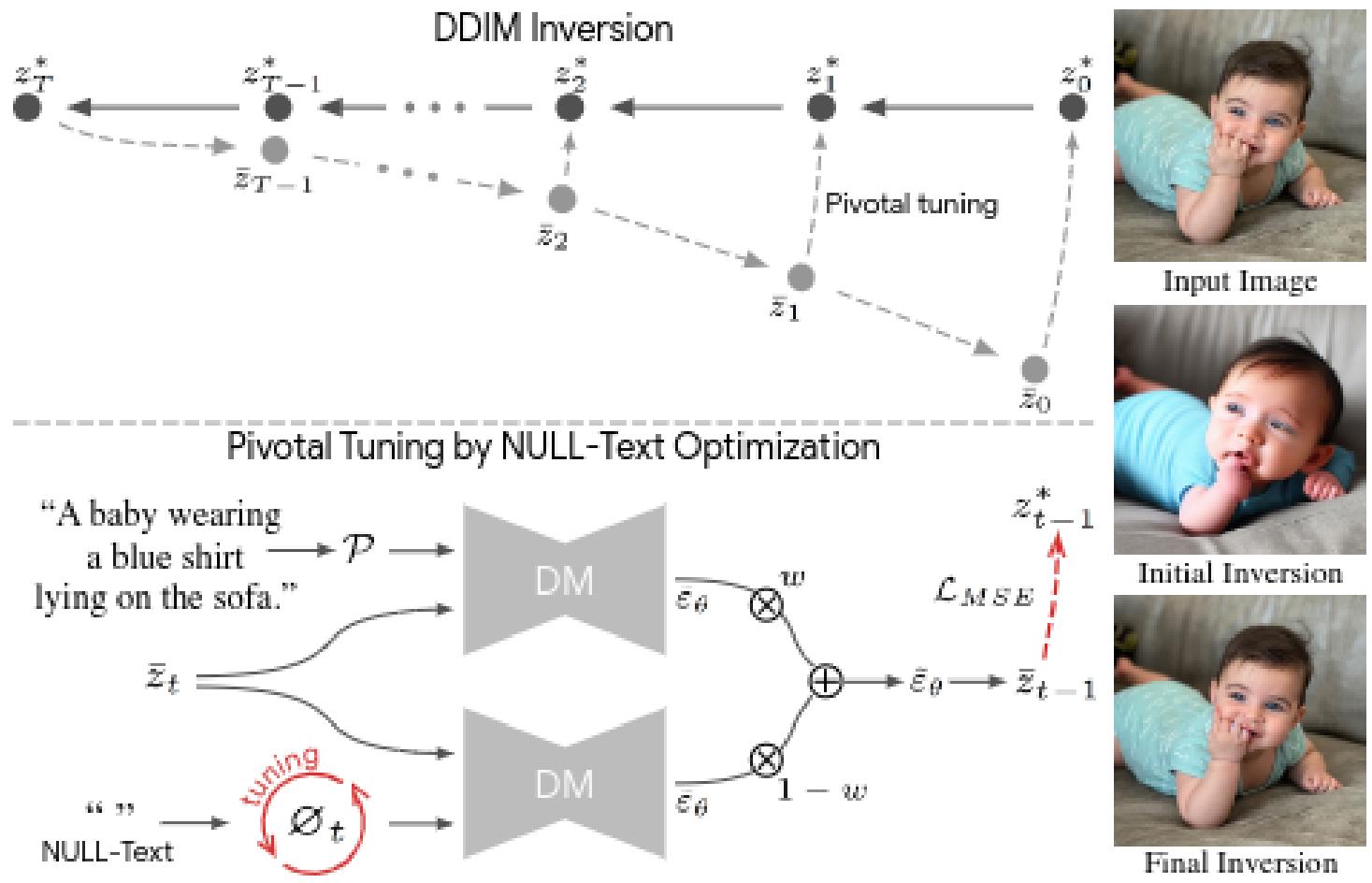


Figure 3. NULL-text Inversion overview. *Top: pivotal inversion.* We first apply an initial DDIM inversion on the input image which estimates a diffusion trajectory $\{z_t^*\}_0^T$. Starting the diffusion process from the last latent z_T^* results in unsatisfying reconstruction as the latent codes become farther away from the original trajectory. We use the initial trajectory as a pivot for our optimization which brings the diffusion backward trajectory $\{\bar{z}_t\}_1^T$ closer to the original image encoding z_0^* . *Bottom: NULL-text optimization for timestamp t .* Recall that classifier-free guidance consists of performing the prediction ϵ_θ twice – using text condition embedding and unconditionally using NULL-text embedding \emptyset (bottom-left). Then, these are extrapolated with guidance scale w (middle). We optimize only the unconditional embeddings \emptyset_t by employing a reconstruction MSE loss (in red) between the predicated latent code z_{t-1} to the pivot z_{t-1}^* .



Algorithm

Algorithm 1: Null-text inversion

- 1 **Input:** A source prompt embedding $\mathcal{C} = \psi(\mathcal{P})$ and input image \mathcal{I} .
 - 2 **Output:** Noise vector z_T and optimized embeddings $\{\emptyset_t\}_{t=1}^T$.
 - 3 Set guidance scale $w = 1$;
 - 4 Compute the intermediate results z_T^*, \dots, z_0^* using DDIM inversion over \mathcal{I} ;
 - 5 Set guidance scale $w = 7.5$;
 - 6 Initialize $\bar{z}_T \leftarrow z_T^*$, $\emptyset_T \leftarrow \psi("")$;
 - 7 **for** $t = T, T - 1, \dots, 1$ **do**
 - 8 **for** $j = 0, \dots, N - 1$ **do**
 - 9 $\emptyset_t \leftarrow \emptyset_t - \eta \nabla_{\emptyset} \|z_{t-1}^* - z_{t-1}(\bar{z}_t, \emptyset_t, \mathcal{C})\|_2^2$;
 - 10 **end**
 - 11 Set $\bar{z}_{t-1} \leftarrow z_{t-1}(\bar{z}_t, \emptyset_t, \mathcal{C})$, $\emptyset_{t-1} \leftarrow \emptyset_t$;
 - 12 **end**
 - 13 **Return** $\bar{z}_T, \{\emptyset_t\}_{t=1}^T$
-



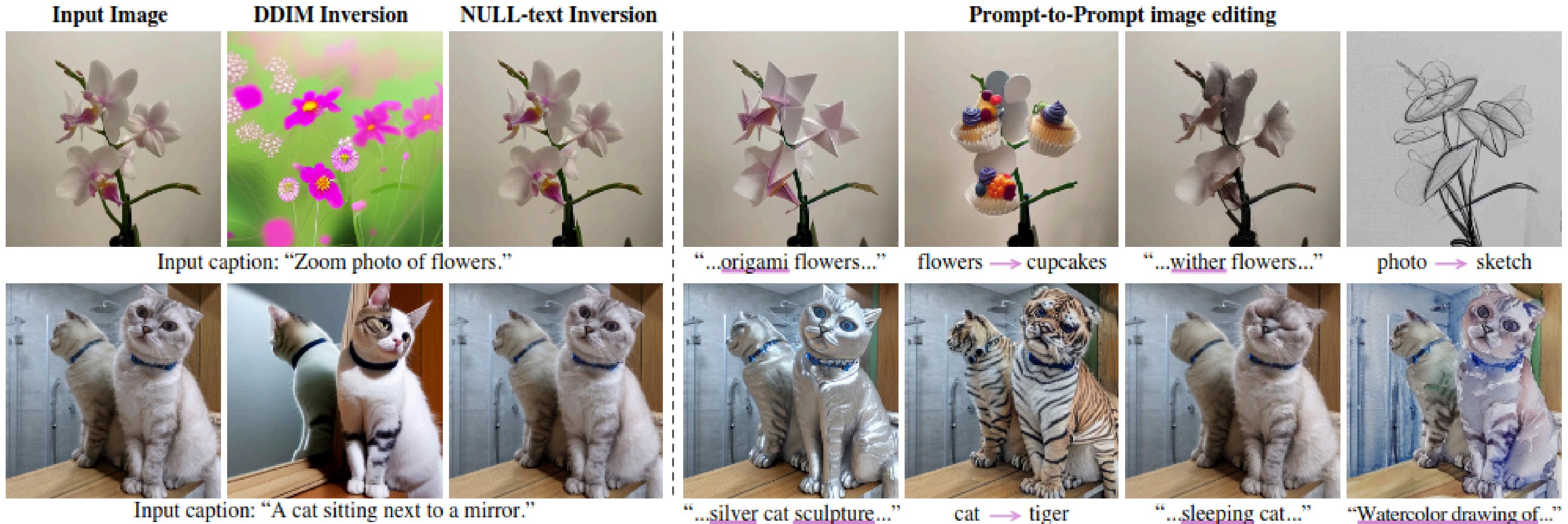


Figure 1. **Null-text inversion for real image editing.** Our method takes as input a real image (leftmost column) and an associated caption. The image is inverted with a DDIM diffusion model to yield a diffusion trajectory (second column to the left). Once inverted, we use the initial trajectory as a pivot for null-text optimization that accurately reconstructs the input image (third column to the left). Then, we can edit the inverted image by modifying only the input caption using the editing technique of Prompt-to-Prompt [18].

Input caption: “A baby wearing a blue shirt lying on the sofa.”



Input caption: “A man in glasses eating a doughnut in the park.”

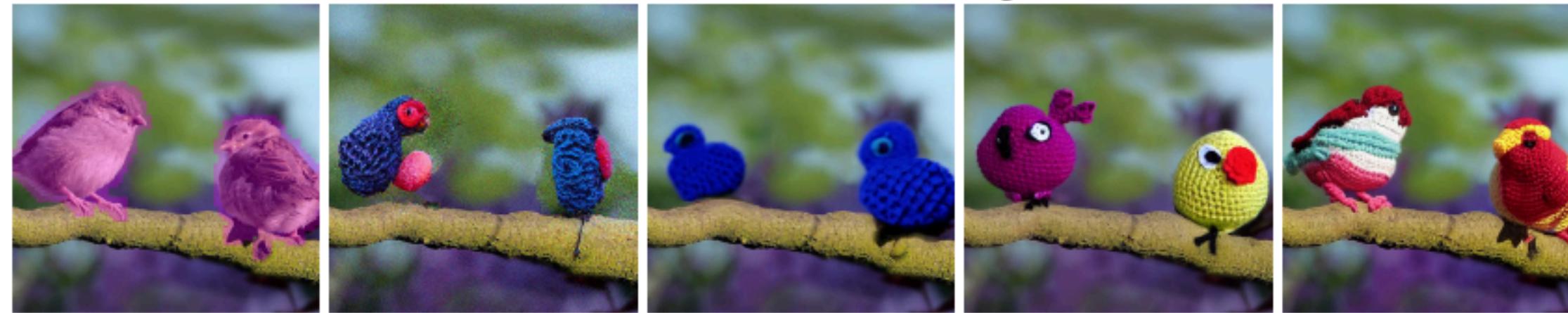


Figure 2. Real image editing using our method. We first apply a single *NULL*-text inversion over the real input image, achieving high-fidelity reconstruction. Then, various Prompt-to-Prompt text-based editing operations are applied. As can be seen, our inversion scheme provides high fidelity while retaining high editability. See additional examples in the supplementary materials.



Figure 6. Comparison. *Text2LIVE* [6] excels at replacing textures locally but struggles to perform more structured editing, such as replacing a kid with a tiger. *VQGAN+CLIP* [13] obtains inferior realism. *SDEdit* [25] fails to faithfully reconstruct the original image, resulting in identity drift when humans are involved. Our method achieves realistic editing of both textures and structured objects while retaining high fidelity to the original image.

Input caption: “Two crochet birds sitting on a branch.”



Input caption: “A basket with apples kittens on a chair.”



Input Image+Mask **Blended-Diffusion**

Glide

SD Inpainting

Ours

Figure 7. Comparison to mask-based methods . As can be seen, mask-based methods do not require inversion as the region outside the mask is kept. However, unlike our approach, such methods often struggle to preserve details that are found inside the masked region. For example, basket size is not preserved.



Evaluation

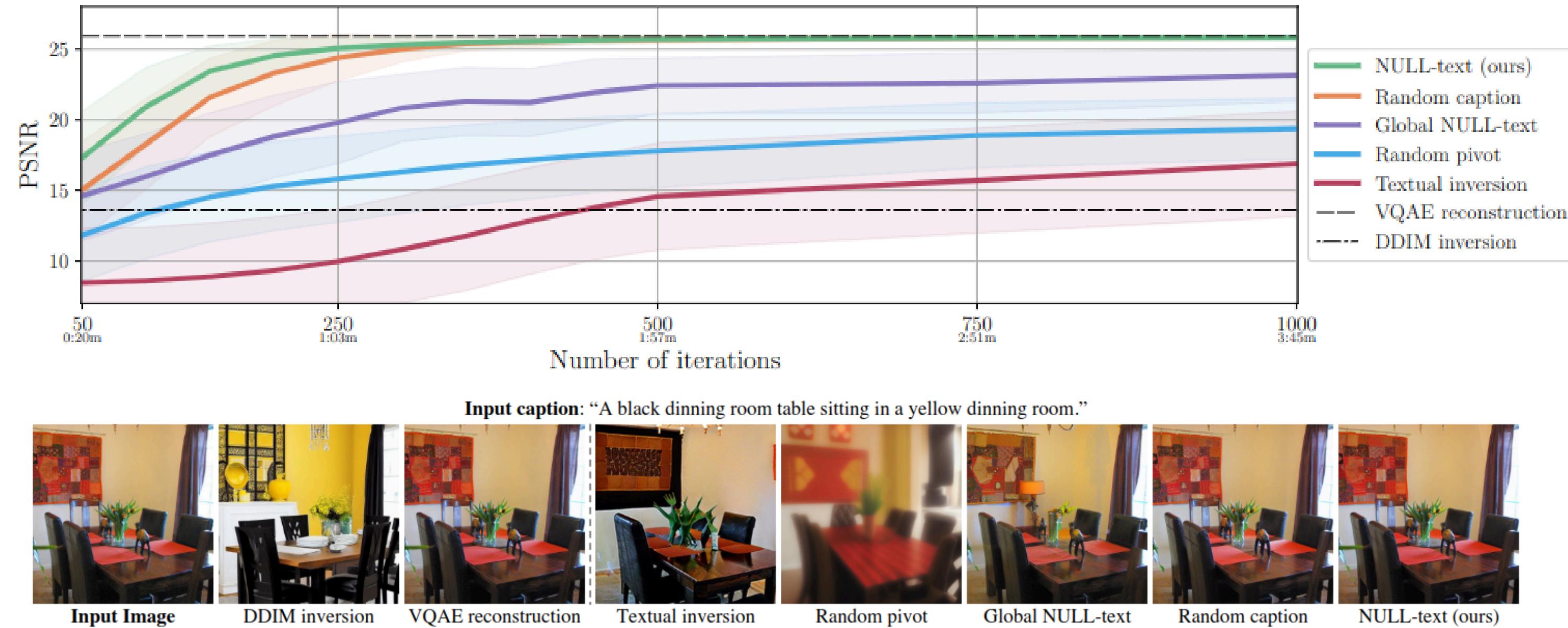


Figure 4. **Ablation Study.** *Top: we compare the performance of our full algorithm (green line) to different variations, evaluating the reconstruction quality by measuring the PSNR score as a function of number optimization iterations and running time in minutes. Bottom: we visually show the inversion results after 200 iterations of our full algorithm (on right) compared to other baselines. Results for all iterations are shown in the supplementary materials.*



Evaluation

Table 2. Quantitative evaluation on ImageNet.

	SDEdit [25]	Text2Live [6]	FlexIT [10]	Ours
LPIPS ↓	31.5	20.0	24.7	29.1
Acc.% ↑	36.5	47.5	51.3	61.1

Table 1. User study results. *The participants were asked to select the best editing result in terms of fidelity to both the input image and the textual edit instruction.*

VQGAN+CLIP	Text2Live	SDEDIT	Ours
3.8%	16.6%	14.5%	65.1%



Prompt-to-Prompt Image Editing with Cross Attention Control

Amir Hertz^{* 1,2}, Ron Mokady^{* 1,2}, Jay Tenenbaum¹, Kfir Aberman¹, Yael Pritch¹, and Daniel Cohen-Or^{* 1,2}

¹ *Google Research*

² *The Blavatnik School of Computer Science, Tel Aviv University*



Prompt-to-Prompt

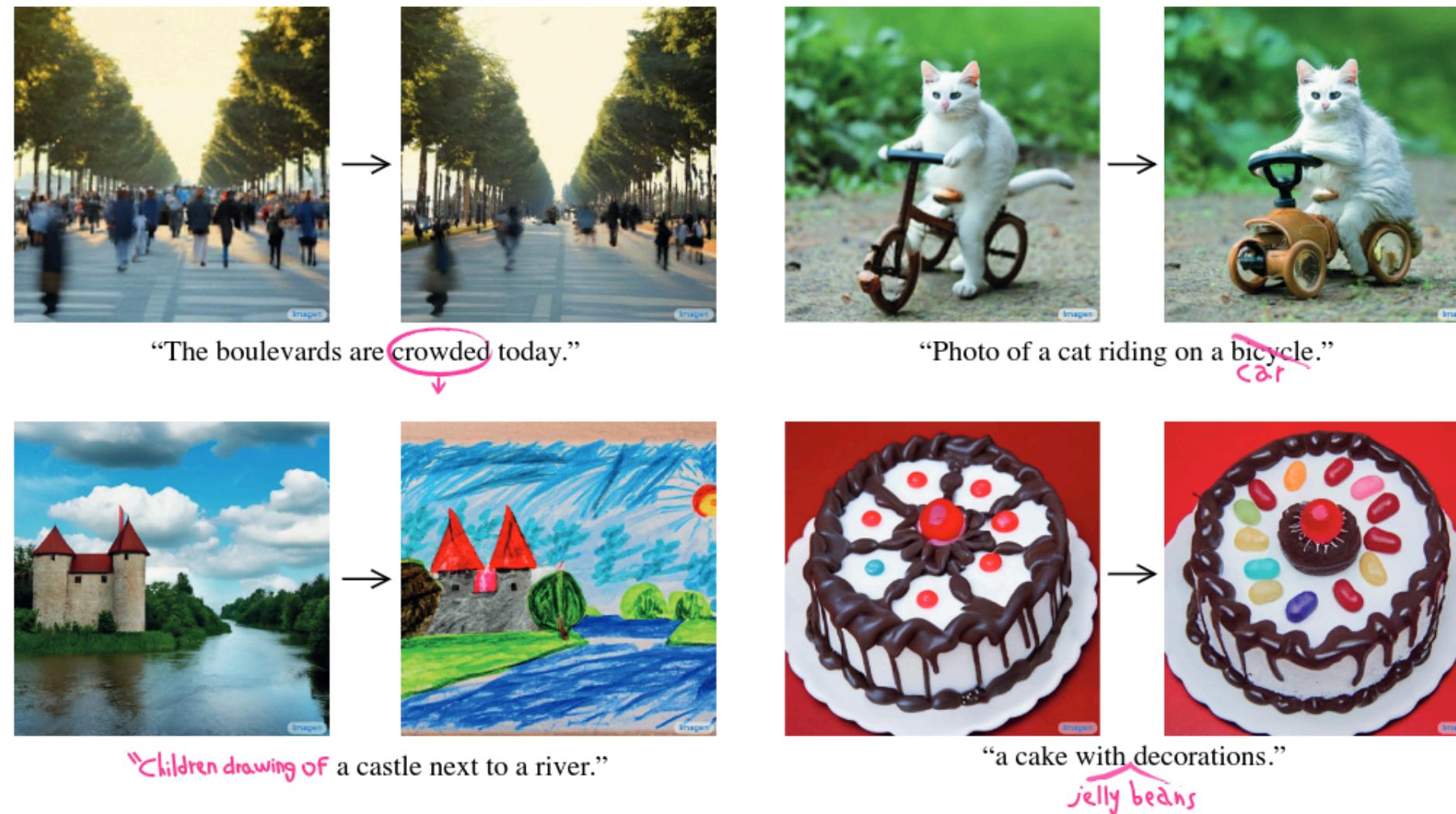


Figure 1: Our method provides variety of *Prompt-to-Prompt* editing capabilities. The user can tune the level of influence of an adjective word (top-left), replace items in the image (top-right), specify a style for an image (bottom-left), or make further refinements over the generated image (bottom-right). The manipulations are infiltrated through the cross-attention mechanism of the diffusion model without the need for any specifications over the image pixel space.



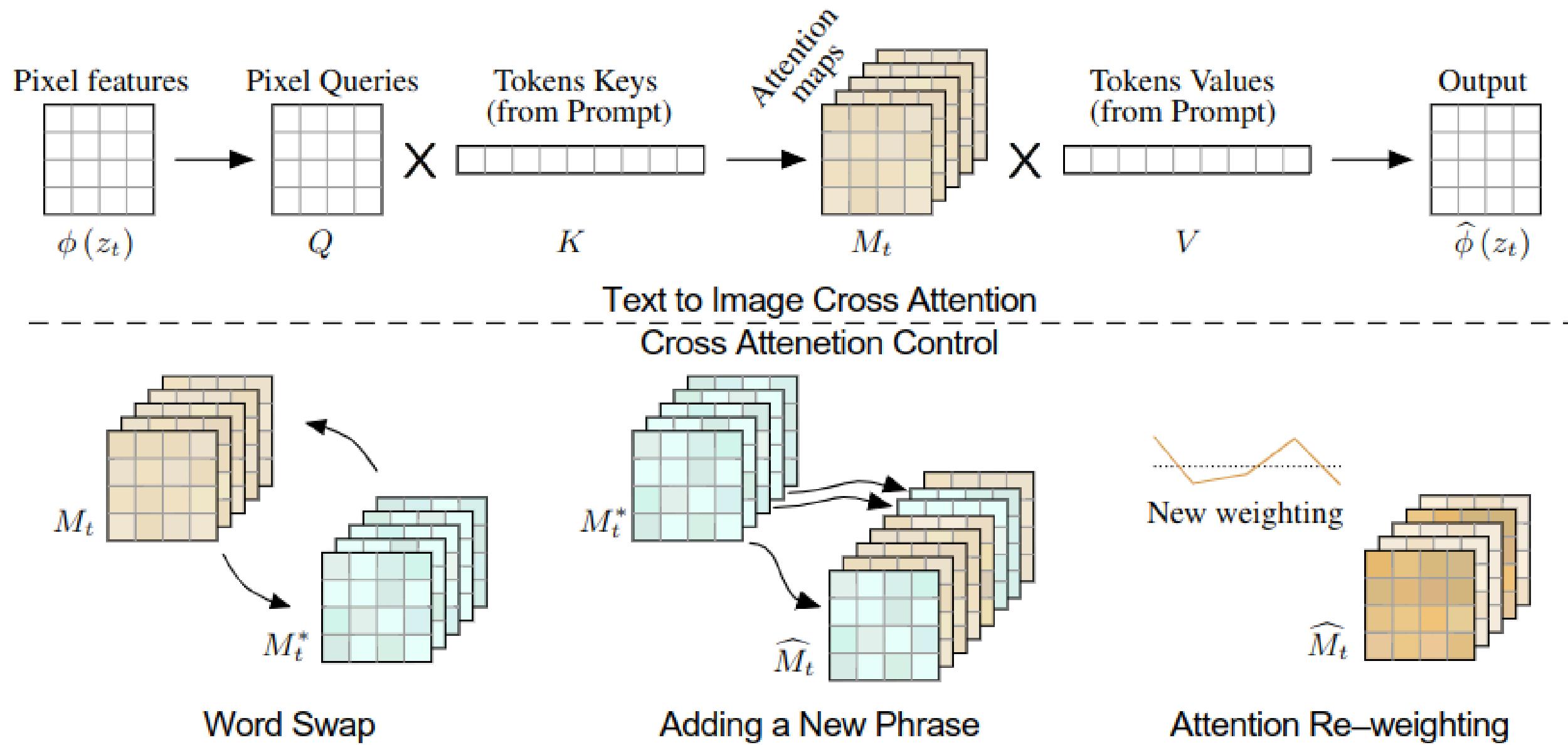


Figure 3: Method overview. Top: visual and textual embedding are fused using cross-attention layers that produce spatial attention maps for each textual token. Bottom: we control the spatial layout and geometry of the generated image using the attention maps of a source image. This enables various editing tasks through editing the textual prompt only. When swapping a word in the prompt, we inject the source image maps M_t , overriding the target image maps M_t^* , to preserve the spatial layout. Where in the case of adding a new phrase, we inject only the maps that correspond to the unchanged part of the prompt. Amplify or attenuate the semantic effect of a word achieved by re-weighting the corresponding attention map.

$$M = \text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right),$$

Algorithm 1: Prompt-to-Prompt image editing

```

1 Input: A source prompt  $\mathcal{P}$ , a target prompt  $\mathcal{P}^*$ , and a random seed  $s$ .
2 Output: A source image  $x_{src}$  and an edited image  $x_{dst}$ .
3  $z_T \sim N(0, I)$  a unit Gaussian random variable with random seed  $s$ ;
4  $z_T^* \leftarrow z_T$ ;
5 for  $t = T, T - 1, \dots, 1$  do
6    $z_{t-1}, M_t \leftarrow DM(z_t, \mathcal{P}, t, s)$ ;
7    $M_t^* \leftarrow DM(z_t^*, \mathcal{P}^*, t, s)$ ;
8    $\widehat{M}_t \leftarrow Edit(M_t, M_t^*, t)$ ;
9    $z_{t-1}^* \leftarrow DM(z_t^*, \mathcal{P}^*, t, s_t) \{ M \leftarrow \widehat{M}_t \}$ ;
10 end
11 Return  $(z_0, z_0^*)$ 

```

word swap

$$Edit(M_t, M_t^*, t) := \begin{cases} M_t^* & \text{if } t < \tau \\ M_t & \text{otherwise.} \end{cases}$$

adding a new phrase

$$(Edit(M_t, M_t^*, t))_{i,j} := \begin{cases} (M_t^*)_{i,j} & \text{if } A(j) = None \\ (M_t)_{i,A(j)} & \text{otherwise.} \end{cases}$$

attention re-weighting

$$(Edit(M_t, M_t^*, t))_{i,j} := \begin{cases} c \cdot (M_t)_{i,j} & \text{if } j = j^* \\ (M_t)_{i,j} & \text{otherwise.} \end{cases}$$

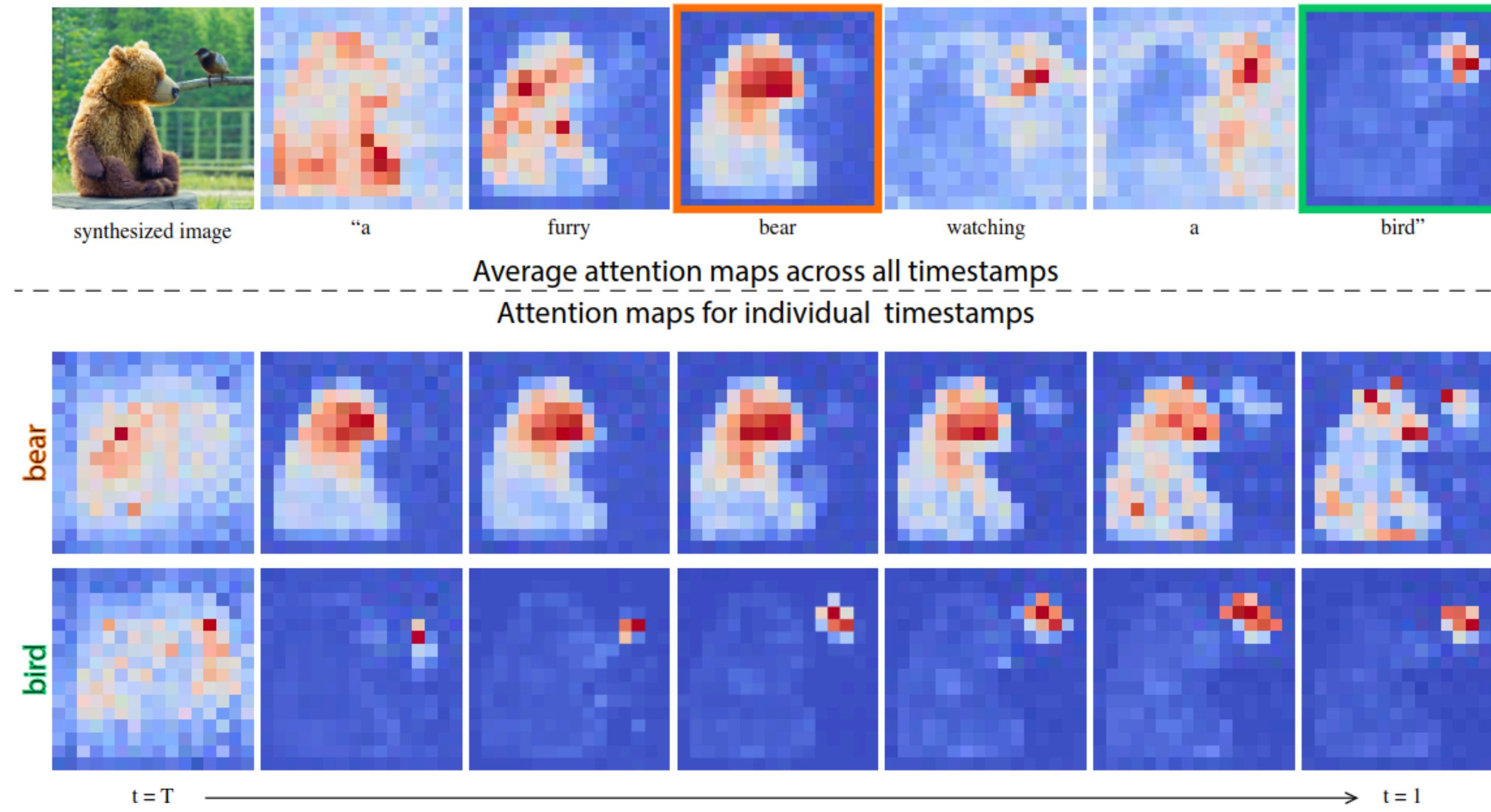
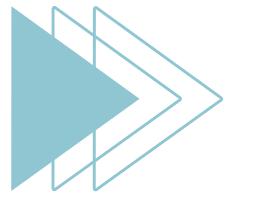


Figure 4: Cross-attention maps of a text-conditioned diffusion image generation. The top row displays the average attention masks for each word in the prompt that synthesized the image on the left. The bottom rows display the attention maps from different diffusion steps with respect to the words “bear” and “bird”.



Questions?

