

---

## Appendix: Energy-Based Cross Attention for Bayesian Context Update in Text-to-Image Diffusion Models

---

388  
389

### 390 A Proof of Theorem 1

391 **Theorem 1.** *For the energy functions*

$$E(\mathbf{Q}; \mathbf{K}) = \frac{\alpha}{2} \text{diag}(\mathbf{K}\mathbf{K}^T) - \sum_{i=1}^N \text{logsumexp}(\mathbf{Q}\mathbf{k}_i^T, \beta) \quad (17)$$

392 *and*

$$E(\mathbf{K}) = \log \sum_{i=1}^N \exp\left(\frac{1}{2} \mathbf{k}_i \mathbf{k}_i^T\right), \quad (18)$$

393 *the gradient of the log posterior is given by:*

$$\nabla_{\mathbf{K}} \log p(\mathbf{K} | \mathbf{Q}) = \text{softmax}_2(\beta \mathbf{K} \mathbf{Q}^T) \mathbf{Q} - \left( \alpha \mathbf{I} + \mathbf{D}\left(\text{softmax}\left(\frac{1}{2} \text{diag}(\mathbf{K}\mathbf{K}^T)\right)\right) \right) \mathbf{K}, \quad (19)$$

394 *Then, by using the chain rule the update rule of context vectors  $\mathbf{C}$  is derived as follows:*

$$\mathbf{C}_{n+1} = \mathbf{C}_n + \gamma \left( \text{softmax}_2(\beta \mathbf{K} \mathbf{Q}^T) \mathbf{Q} - \left( \alpha \mathbf{I} + \mathbf{D}\left(\text{softmax}\left(\frac{1}{2} \text{diag}(\mathbf{K}\mathbf{K}^T)\right)\right) \right) \mathbf{K} \right) \mathbf{W}_K^T, \quad (20)$$

395 *where  $\gamma > 0$  is a step size, and  $\mathbf{D}(\cdot)$  is a vector-to-diagonal-matrix operator.*

396 *Proof.* Based on the Bayes' theorem, the gradient of the log posterior is derived as:

$$\nabla_{\mathbf{K}} \log p(\mathbf{K} | \mathbf{Q}) = -(\nabla_{\mathbf{K}} E(\mathbf{Q}; \mathbf{K}) + \nabla_{\mathbf{K}} E(\mathbf{K})). \quad (21)$$

397 First, with definition (17),

$$\nabla_{\mathbf{K}} E(\mathbf{Q}; \mathbf{K}) = \alpha \mathbf{K} - \nabla_{\mathbf{K}} \sum_{i=1}^N \text{logsumexp}(\mathbf{Q}\mathbf{k}_i^T, \beta), \quad (22)$$

398 where  $\forall i \in \{1, \dots, N\}$ ,

$$\begin{aligned} \nabla_{\mathbf{k}_i} \sum_{i=1}^N \text{logsumexp}(\mathbf{Q}\mathbf{k}_i^T, \beta) &= \frac{1}{\beta} \nabla_{\mathbf{k}_i} \log \sum_{j=1}^{P_i^2} \exp(\beta \mathbf{q}_j \mathbf{k}_i^T) \\ &= \sum_{j=1}^{P_i^2} \frac{\exp(\beta \mathbf{q}_j \mathbf{k}_i^T)}{\sum_{n=1}^{P_i^2} \exp(\beta \mathbf{q}_n \mathbf{k}_i^T)} \mathbf{q}_j \\ &= \text{softmax}(\mathbf{Q}\mathbf{k}_i^T)^T \mathbf{Q}. \end{aligned} \quad (23)$$

399 Then, by considering that  $\mathbf{k}_i$  is a  $i$ -th row vector of  $\mathbf{K}$ ,

$$\begin{aligned} \nabla_{\mathbf{K}} \sum_{i=1}^N \text{logsumexp}(\mathbf{Q}\mathbf{k}_i^T, \beta) &= (\text{softmax}_1(\beta \mathbf{Q}\mathbf{K}^T))^T \mathbf{Q} \\ &= \text{softmax}_2(\beta \mathbf{K}\mathbf{Q}^T) \mathbf{Q}, \end{aligned} \quad (24)$$

400 where the last equality holds due to the definition of  $\text{softmax}_1$  in Section 2.2.

401 Second, with definition (18),  $\nabla_{\mathbf{K}} \mathbb{E}(\mathbf{K}) = \nabla_{\mathbf{K}} \log \sum_{i=1}^N \exp(\frac{1}{2} \mathbf{k}_i \mathbf{k}_i^T)$ , where

$$\begin{aligned} \nabla_{\mathbf{k}_i} \log \sum_{i=1}^N \exp(\frac{1}{2} \mathbf{k}_i \mathbf{k}_i^T) &= \frac{\exp(\frac{1}{2} \mathbf{k}_i \mathbf{k}_i^T)}{\sum_{j=1}^N \exp(\frac{1}{2} \mathbf{k}_j \mathbf{k}_j^T)} \mathbf{k}_i \\ &= \text{softmax} \left( \frac{1}{2} \text{diag}(\mathbf{K} \mathbf{K}^T) \right)_i \mathbf{k}_i, \end{aligned} \quad (25)$$

402 where  $\text{softmax}(\cdot)_i$  denotes  $i$ -th value of a softmax vector. Then,

$$\nabla_{\mathbf{K}} \log \sum_{i=1}^N \exp(\frac{1}{2} \mathbf{k}_i \mathbf{k}_i^T) = \mathbf{D} \left( \text{softmax} \left( \frac{1}{2} \text{diag}(\mathbf{K} \mathbf{K}^T) \right) \right) \mathbf{K}, \quad (26)$$

403 where  $\mathbf{D}(\cdot)$  is a vector-to-diagonal-matrix operator that takes  $N$ -dimensional softmax vector as an  
 404 input and returns a  $N \times N$  diagonal matrix with softmax values as main diagonal entries. Then, By  
 405 combining (22), (24) and (26), one can finally obtain:

$$\nabla_{\mathbf{K}} \log p(\mathbf{K} | \mathbf{Q}) = \text{softmax}_2 (\beta \mathbf{K} \mathbf{Q}^T) \mathbf{Q} - \left( \alpha \mathbf{I} + \mathbf{D} \left( \text{softmax} \left( \frac{1}{2} \text{diag}(\mathbf{K} \mathbf{K}^T) \right) \right) \right) \mathbf{K}. \quad (27)$$

406 By using the chain rule with  $\mathbf{K} = \mathbf{C} \mathbf{W}_K$ , the update rule of context vectors  $\mathbf{C}$  is derived as in  
 407 (20).  $\square$

408 We introduce vector-to-matrix operator  $\mathbf{D}(\cdot)$  to avoid confusion and fix the typo in the main paper.

## 409 B Pseudo-code for BCU and CACAO

410 This section provides the description of the pseudocode for the proposed Bayesian Context Update  
 411 (BCU) and Compositional Averaging of Cross-Attention Output (CACAO). Algorithm 1 outlines  
 412 the cascaded context propagation across cross-attention layers within the UNet model during the  
 413 sampling step  $t$ . Note that the context is reinitialized at the beginning of each sampling step. On the  
 414 other hand, Algorithm 2 details the BCU implemented in each cross-attention layer. Remark that  
 415  $\mathbf{D}$  in line 5 denotes vector-to-diagonal-matrix operator. Specifically, the proposed BCU provides a  
 416 significant computational efficiency by reusing the similarity  $\mathbf{Q} \mathbf{K}^T$ , which requires computational  
 417 cost  $\mathcal{O}(N^2)$ , to compute  $\nabla_{\mathbf{K}} E(\mathbf{Q}; \mathbf{K})$ . Consequently, there is only a small amount of additional  
 418 computational overhead associated with the proposed BCU.

---

### Algorithm 1 Context cascade at sampling step $t$

---

**Require:**  $\mathbf{Q}_t, \mathbf{C}_{clip}, \text{UNet}$   
 1:  $\mathbf{C}_t \leftarrow \mathbf{C}_{clip}$  // Re-initialize  
 2: **for** layer in UNet **do**  
 3:     **if** layer is CrossAttention **then**  
 4:          $\mathbf{Q}_t, \mathbf{C}_t \leftarrow \text{layer}(\mathbf{Q}_t, \mathbf{C}_t)$  // Algorithm 2  
 5:     **else**  
 6:          $\mathbf{Q}_t \leftarrow \text{layer}(\mathbf{Q}_t)$   
 7:     **end if**  
 8: **end for**  
 9:  $\mathbf{Q}_{t+1} \leftarrow \mathbf{Q}_t$   
 10: **return**  $\mathbf{Q}_{t+1}$

---

---

**Algorithm 2** Bayesian Context Update (BCU)

---

**Require:**  $Q, C, W_q, W_k, W_v, \alpha, \beta, \gamma_{\text{attn}}, \gamma_{\text{reg}}$   
1:  $Q, K, V \leftarrow QW_q, CW_k, CW_v$   
2:  $S = QK^T$   
3:  $Q \leftarrow \text{softmax}_2(\beta S)V$   
4:  $\nabla_K E(Q; K) = \text{softmax}_2(\beta S^T)Q$   
5:  $\nabla_K E(K) = -(\alpha I + D(\text{softmax}(\frac{1}{2} \text{diag}(KK^T))))K$   
6:  $\Delta C = (\gamma_{\text{attn}} \nabla_K E(Q; K) + \gamma_{\text{reg}} \nabla_K E(K))W_k^T$   
7:  $C \leftarrow C + \Delta C$   
8: **return**  $Q, C$

---

419 Algorithm 3 outlines the pseudocode for the CACAO implemented for  $M$  given contexts. For the  
420 simplicity, we exclude the BCU from the algorithm. Nonetheless, the BCU and the CACAO could be  
421 leveraged together.

---

**Algorithm 3** Compositional Averaging of Cross-Attention Output (CACAO)

---

**Require:**  $Q, C = \{C_1, \dots, C_M\}, W_q, W_k, W_v, \alpha_s, \beta$   
1:  $Q \leftarrow QW_q$   
2: **for**  $s$  in  $[1, \dots, M]$  **do**  
3:      $K_s, V_s \leftarrow C_s W_k, C_s W_v$   
4:      $S_s = QK_s^T$   
5: **end for**  
6:  $Q \leftarrow \frac{1}{M} \sum_{s=1}^M \alpha_s \text{softmax}_2(\beta S_s)V_s$   
7: **return**  $Q$

---

## 422 C Experimental setups

423 In this section, we describe detailed experimental setups for three applications including baseline  
424 method, hyper-parameter of the proposed method, and dataset if it is the case. Code: <https://github.com/EnergyAttention/Energy-Based-CrossAttention>.  
425

### 426 C.1 Common experimental setup

427 We mainly leverage pre-trained Stable Diffusion v1-5 (except Table 1: v1-4) which is provided by  
428 *diffusers*, a Python library that offers various Stable Diffusion pipelines with pre-trained models.  
429 All images are sampled for 50 steps via PNDM sampler [20] using NVIDIA RTX 2080Ti. In every  
430 experiment, we set the parameter  $\alpha$  in Equation (18) to zero, focusing solely on controlling the values  
431 of  $\gamma_{\text{attn}}$  and  $\gamma_{\text{reg}}$ . BCU is applied to every task, and CACAO is additionally employed in C.4.

432 **Different learning rate for each token** It is worth noting that the  $\gamma_{\text{attn}}$  and  $\gamma_{\text{reg}}$  could be expressed  
433 as vectors. In other words, if the context  $C \in \mathbb{R}^{N \times d_c}$  is given,  $\gamma_{\text{attn}}$  and  $\gamma_{\text{reg}}$  are  $N$ -dimensional  
434 vectors. Hence, we have the flexibility to adjust the learning rate  $\gamma_{\{\cdot\}}$ , allowing us to increase or  
435 decrease the impact of certain tokens based on the user’s intent. Unless otherwise noted,  $\gamma_{\text{attn}}$  and  
436  $\gamma_{\text{reg}}$  is set to a constant for each text token.

437 **Learning rate scheduling** Since the proposed BCU is leveraged for the diffusion model, one can  
438 readily introduce scheduling strategies for  $\gamma_{\text{attn}}$  and  $\gamma_{\text{reg}}$  along the sampling step  $t$ . We implement  
439 multiple variants such as ‘constant’, ‘step’, and ‘exponential decay’ as follows.

$$\begin{aligned} \text{[constant]} \quad & \gamma(t) = \gamma_0 \\ \text{[step]} \quad & \gamma(t) = \gamma_0 \cdot \text{ReLU}(t - \tau) \\ \text{[exp-decay]} \quad & \gamma(t) = \gamma_0 \cdot \lambda^t \end{aligned} \tag{28}$$

440 where  $\gamma_0$  is the initial value,  $\text{ReLU}(x) = 0$  if  $x \leq 0$ , otherwise 1,  $\tau$  denotes the temporal threshold,  
441 and  $\lambda$  denotes the decay ratio. Unless stated otherwise, the scheduling strategy is set to the ‘constant’.

## 442 C.2 Multi-concept image generation

443 We compared the performance of the proposed method with Structured Diffusion [11] which does not  
444 require additional training as our method. We leveraged the open-sourced official implementation <sup>1</sup>.

445 For the proposed method, we set the  $\gamma_{attn}$  and  $\gamma_{reg}$  differently for each sample within [1e-2, 1.5e-2,  
446 2e-2]. As shown in the following ablation studies E, large  $\gamma_{attn}$  tends to generate saturated images  
447 while large  $\gamma_{reg}$  results in mixed/vanished contents.

448 We found that using different learning rates for each context token is useful for multi-concept  
449 generation, especially when a single concept tends to dominate with a constant learning rate. For  
450 example, given the main prompt "A cat wearing a shirt", we set the  $\gamma_{attn}$  for the "shirt" to  
451 3e-2, while  $\gamma_{attn}$  is set to 1.5e-2 for other tokens. We have observed that doubling the  $\gamma_{attn}$  for a text  
452 token to be emphasized is sufficient to achieve balanced multi-concept image generation for most  
453 cases.

## 454 C.3 Text-guided image inpainting

455 Additionally, we conducted a performance comparison between our proposed method and two  
456 alternative approaches: (a) Stable Inpaint<sup>2</sup>, which fine-tunes the weights of Stable Diffusion through  
457 inpainting training, and (b) Stable Repaint<sup>3</sup>, which leverages the work of Lugmayr et al. [22] on the  
458 latent space of Stable Diffusion for the inpainting task. In the case of Stable Repaint, the mask is  
459 downsized and transferred into the latent space. We applied the Bayesian Context Update (BCU)  
460 technique to both methods, resulting in improved results compared to their respective baselines.

461 **Masked BCU.** To further enhance the performance for the inpainting task, we introduce the concept  
462 of masked Bayesian Context Update (masked BCU). Specifically, let  $M \in \mathbb{R}^{P_l^2 \times P_l^2}$  represent a  
463 diagonal matrix where the main diagonal values are derived from the downsampled inpainting mask  
464 for the  $l$ -th cross-attention layer, with an output spatial size of  $P_l^2$ . In Equation (29), we modify the  
465 attention term (12) by incorporating the downsampled mask, effectively covering the query matrix as  
466 follows:

$$C_{n+1} = C_n + \gamma \left( \text{softmax}_2(\beta K Q^T) M Q - \left( \alpha I + D \left( \text{softmax} \left( \frac{1}{2} \text{diag}(K K^T) \right) \right) \right) K \right) W_K^T. \quad (29)$$

467 As evident in Equation (12), the attention term updates the context vectors, aligning  $k_i$  towards  
468  $q_j, j = 1, \dots, P_l^2$ , while considering the alignment strength between each  $q_j$  and  $k_i$ . However,  
469 in the inpainting task, we have prior knowledge that the context vectors should be most aligned  
470 with the semantically relevant masked regions. Therefore, we mask out unrelated background  
471 spatial representations, allowing for the context vectors to be updated with a specific focus on the  
472 masked regions. This approach facilitates the incorporation of semantic information encoded by  $k_i$   
473 specifically into the spatial mask regions.

474 In our proposed method, we set different values for  $\gamma_{attn}$  and  $\gamma_{reg}$  for each sample, selected from the  
475 set [1e-2, 1.5e-2, 2e-2, 2.5e-2], to account for variations in the input samples.

## 476 C.4 Image editing via compositional generation

477 We present empirical evidence demonstrating the effectiveness of our energy-based framework for  
478 compositional synthetic and real-image editing. The Bayesian Context Update (BCU) technique can  
479 be readily applied to both the main context vector ( $C_1$  in Section 3.2,  $s = 1$ ) and editorial context  
480 vectors ( $C_{s>1}$ ). Each BCU operation influences the attention maps used in Compositional Averaging  
481 of Cross-Attention Output (CACAO), enhancing the conveyance of semantic information associated  
482 with each context. Note that  $\alpha_s$  in (16) represents the degree of influence of the  $s$ -th concept in the  
483 composition. In practice, we fix  $\alpha_1 = 1$  for the main context, while  $\alpha_{s>1}$  is tuned within the range of  
484 (0.5, 1.0).

<sup>1</sup><https://github.com/weixi-feng/Structured-Diffusion-Guidance>

<sup>2</sup><https://huggingface.co/runwayml/stable-diffusion-inpainting>

<sup>3</sup><https://github.com/huggingface/diffusers/tree/main/examples/community#stable-diffusion-repaint>

485 Let  $\gamma_{attn,s}$  and  $\gamma_{reg,s}$  denote the step sizes for BCU of the  $s$ -th context vector. If the editing process  
486 involves changing the identity of the original image (e.g., transforming a "cat" into a "dog"), we set  
487 both  $\gamma_{attn,1}$  and  $\gamma_{reg,1}$  to zero. Otherwise, if the editing maintains the original identity, we choose  
488 values for  $\gamma_{attn,1}$  and  $\gamma_{reg,1}$  from the range of (5e-4, 1e-3), similar to  $\gamma_{attn,(s>1)}$  and  $\gamma_{reg,(s>1)}$ . All  
489 hyperparameters, including  $\alpha_s$  and  $\gamma_s$ , are fixed during the quantitative evaluation process (more  
490 details in Section D and Table 2).

491 To ensure consistent results, we maintained a fixed random seed for both real and synthetic image  
492 editing. For real image editing, we employed null-text pivotal inversion [24] to obtain the initial noise  
493 vector.

494 During the reverse diffusion process in Sections C.2 and C.3, we kept  $\gamma$  fixed as a constant value.  
495 However, for compositional generation, we utilized step scheduling (Equation 28) for  $\gamma_s$  and  $\alpha_s$ .  
496 After converting the initial noise vector for real images or using a fixed random seed for synthetic  
497 images, BCU and CACAO are applied after a threshold time  $\tau_s > 0$  for the  $s$ -th editorial context.  
498 This scheduling strategy helps to preserve the overall structure of generated images during the  
499 editing process. In our observations, a value of  $\tau_s \in [10, 25]$  generally produces satisfactory results,  
500 considering a total number of reverse steps set to 50. However, one can increase or decrease  $\tau_s$  for  
501 more aggressive or conservative editing, respectively.

502 The exemplary real images presented in Figures 5 and 6 of the main paper were sampled from datasets  
503 such as FFHQ [16], AFHQ [5], and ImageNet [7]. For a detailed quantitative analysis, please refer to  
504 Section D.

## 505 D Quantitative Comparison

506 In this section, we conducted a comparative analysis of the proposed framework against several  
507 state-of-the-art diffusion-based image editing methods [23, 12, 24, 27], following the experimental  
508 setup of [27]. To ensure a fair comparison, all methods utilize the pre-trained Stable Diffusion v1-4,  
509 employ the PNDM sampler with an equal number of sampling steps, and adopt the same classifier-free  
510 guidance scale.

### 511 D.1 Baseline Methods

512 In addition to the Plug-and-Play method discussed in the main paper, we include the following  
513 baselines for comprehensive quantitative comparison:

514 **SDEdit [23] + word swap.** This method introduces the Gaussian noise of an intermediate timestep  
515 and progressively denoises images using a new textual prompt, where the source word (e.g., Cat) is  
516 replaced with the target word (e.g., Dog).

517 **Prompt-to-prompt (P2P) [12].** P2P edits generated images by leveraging explicit attention maps  
518 from a source image. The source attention maps  $M_t$  are used to inject, re-weight, or override the  
519 target maps based on the desired editing operation. These original maps act as hard constraints for  
520 the edited images.

521 **DDIM + word swap [24].** This method applies null-text inversion to real input images, achieving  
522 high-fidelity reconstruction. DDIM sampling is then performed using inverted noise vectors and an  
523 edited prompt generated by swapping the source word with the target.

524 **pix2pix-zero [27].** pix2pix-zero first derives a text embedding direction vector  $\Delta_{c_{edit}}$  from the source  
525 to the target by using a large bank of diverse sentences generated from a state-of-the-art sentence  
526 generator, such as GPT-3 [2]. Inverted noise vectors are denoised with the edited text embedding,  
527  $c + \Delta_{c_{edit}}$ , and cross-attention guidance to preserve consensus.

### 528 D.2 Dataset

529 For our quantitative evaluations, we focus on three image-to-image translation tasks: (1) translating  
530 cats to dogs (cat  $\rightarrow$  dog), (2) translating horses to zebras (horse  $\rightarrow$  zebra), and (3) adding glasses  
531 to cat input images (cat  $\rightarrow$  cat with glasses). Following the data collection protocol of [27], we  
532 retrieve 250 relevant cat images and 213 horse images from the LAION 5B dataset [33] using CLIP

533 embeddings of the source text description. We select images with a high CLIP similarity to the source  
 534 word for each task.

### 535 D.3 Metrics

536 Motivated by [38, 27], we measure CLIP Accuracy and DINO-ViT structure distance. Specifically,  
 537 (a) CLIP Acc represents whether the targeted semantic contents are well reflected in the generated  
 538 images. It calculates the percentage of instances where the edited image has a higher similarity to  
 539 the target text, as measured by CLIP, than to the original source text [27]. On the other hand, (b)  
 540 structure distance [38, 37] measures whether the overall structure of the input image is well preserved.  
 541 It is defined as the difference in self-similarity of the keys extracted from the attention module at the  
 542 deepest DINO-ViT [3] layer.

### 543 D.4 Details

544 The main context vector  $C_{main}$  is encoded given a main prompt automatically generated by BLIP  
 545 [19]. In addition, the editorial context vectors  $C_{src}$  and  $C_{tgt}$  are encoded given the text descriptions  
 546 of the source and target concept, i.e. source and target prompt. For example, for a cat  $\rightarrow$  dog task (cat  
 547  $\rightarrow$  cat w/ glasses), the source prompt is "cat" ("cat wearing glasses"), and the target prompt  
 548 is "dog" ("without glasses"). Then we apply BCU and CACAO based on the obtained context  
 549 vectors. Please refer to Table 2 for the hyperparameter configurations.

### 550 D.5 Results

551 Table 1 shows that the proposed energy-based framework gets a high CLIP-Acc while having  
 552 low Structure Dist. It implies that the proposed framework can perform the best edit while still  
 553 retaining the structure of the original input image. This is a remarkable result considering that the  
 554 proposed framework is not specially designed for the real-image editing task. Moreover, the proposed  
 555 framework does not rely on the large bank of prompts and editing vector  $\Delta_{C_{edit}}$  [27] which can be  
 556 easily incorporated into our method.

557 While DDIM + word swap records remarkably high CLIP-Acc in horse  $\rightarrow$  zebra task, Figure 7 and  
 558 12 show that such improvements are based on unintended changes in the overall structure. Table 2  
 559 summarizes the hyperparameter settings for each task. Examples of results are presented in Figure 13  
 560 and 12.

Table 1: Comparison to state-of-the-art diffusion-based editing methods. Dist for DINO-ViT Structure distance. Baseline results are from [27].

Method	(a) Cat $\rightarrow$ Dog		(b) Horse $\rightarrow$ Zebra		(c) Cat $\rightarrow$ Cat w/ glasses	
	CLIP Acc ( $\uparrow$ )	Dist ( $\downarrow$ )	CLIP Acc ( $\uparrow$ )	Dist ( $\downarrow$ )	CLIP Acc ( $\uparrow$ )	Dist ( $\downarrow$ )
SDEdit [23] + word swap	71.2%	0.081	92.2%	0.105	34.0%	0.082
DDIM + word swap	72.0%	0.087	<b>94.0%</b>	0.123	37.6%	0.085
prompt-to-prompt [12]	66.0%	0.080	18.4%	0.095	69.6%	0.081
pix2pix-zero [27]	92.4%	0.044	75.2%	0.066	71.2%	<b>0.028</b>
Stable Diffusion + ours	<b>93.7%</b>	<b>0.040</b>	90.4%	<b>0.061</b>	<b>81.1%</b>	0.052

Table 2: Hyperparameter configurations for each editing task. Each task index comes from Table 1.  $\gamma_{attn,main} = 0$  and  $\gamma_{reg,main} = 0$  as mentioned in section C.4. Note that  $\alpha_{src} < 0$  for the concept negation (related ablation study in Figure 9).  $\tau_s$  denotes the warm-up period for step scheduling in (28) and Section C.4.

Task	$\alpha_{src}$	$\alpha_{tgt}$	$\gamma_{\cdot,main}$	$\gamma_{attn,src}$	$\gamma_{reg,src}$	$\gamma_{attn,tgt}$	$\gamma_{reg,tgt}$	$\tau_s$
(a)	-0.65	0.75	0	5e-4	5e-4	6e-4	6e-4	25
(b)	-0.5	0.6	0	4e-4	4e-4	5e-4	5e-4	15
(c)	-0.6	0.7	0	1e-3	1e-3	1e-3	1e-3	17

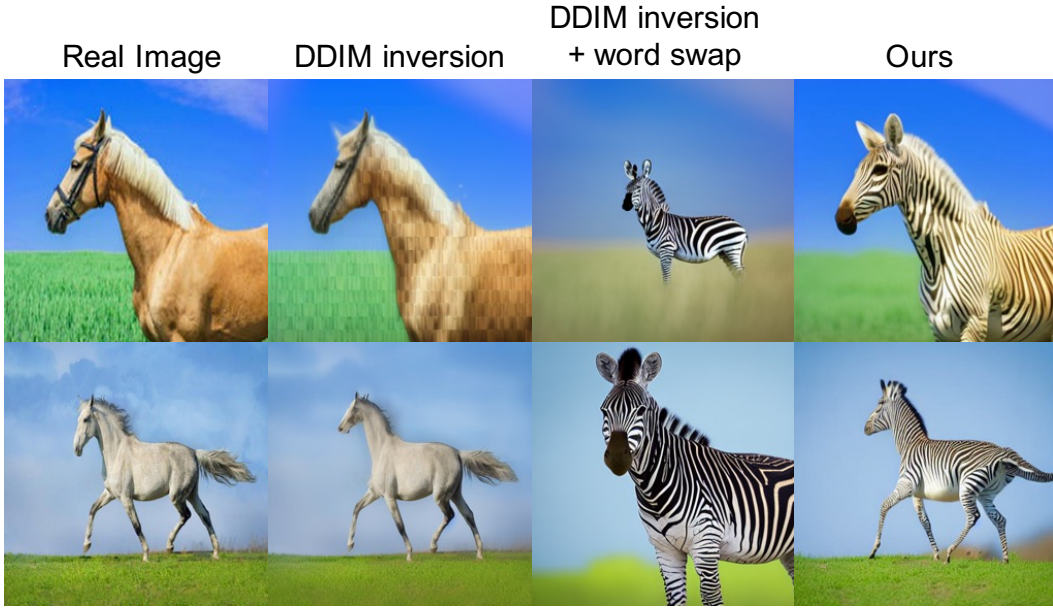


Figure 7: Image editing comparison with DDIM-inversion. Generated samples by DDIM-inversion with word swap readily deviate the original data contents, while the proposed method avoids undesired changes.

561 **E Ablation study and more results**

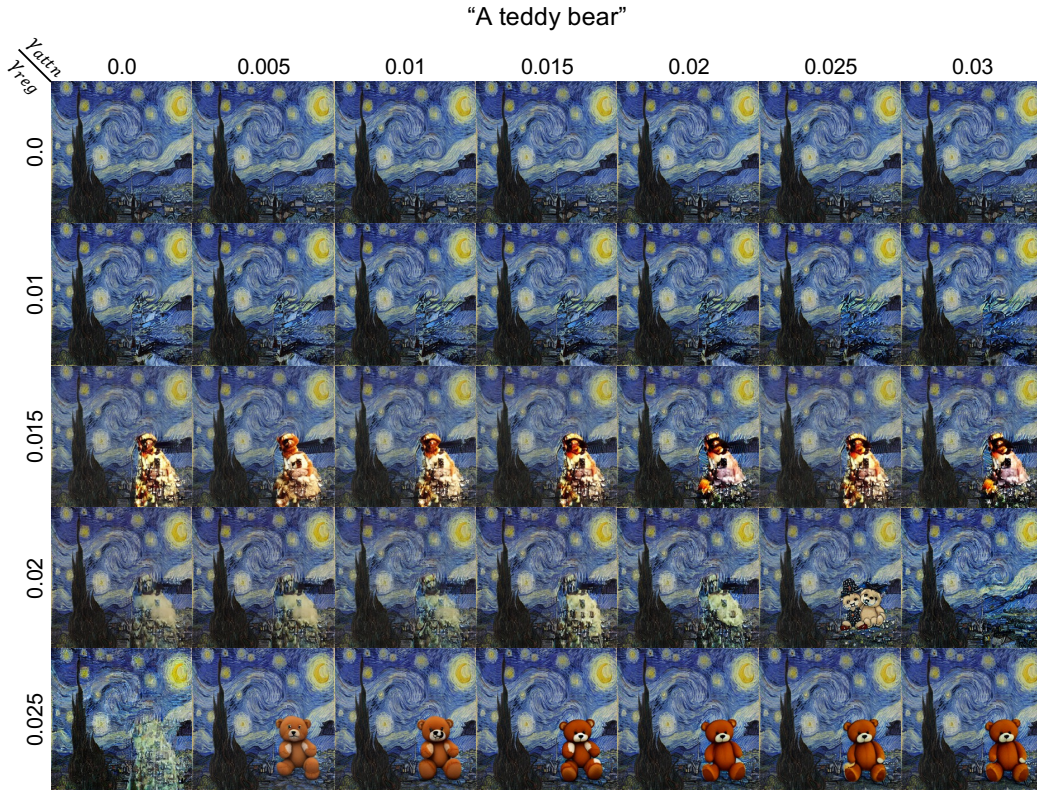


Figure 8: Ablation results for  $\gamma_{attn}$  and  $\gamma_{reg}$ . All samples are generated from the same random noise.



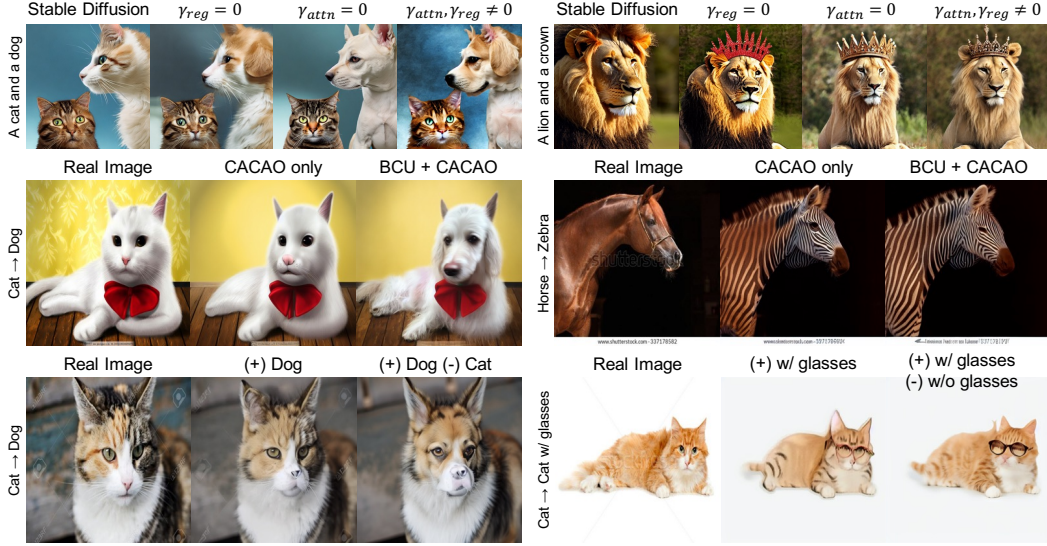


Figure 9: Ablation study results. The first row shows multi-concept generation examples with varying  $\gamma_{attn}$  and  $\gamma_{reg}$ , while the second row shows real image editing examples with varying the usage of BCU and CACAO. The last row shows the effect of negative prompt for the image editing application.

562 **Attention and regularization terms.** To access the degree of performance improvement attained  
563 by the proposed BCU, we conducted an ablation study for the attention and the regularization terms  
564 by regulating  $\gamma_{attn}$  and  $\gamma_{reg}$  for the text-guided image inpainting (Figure 8) and the multi-concept  
565 image generation (Figure 9). From the Figure 8, we can observe that the desired content is generated  
566 when proper range of  $\gamma_{attn}$  and  $\gamma_{reg}$  are given. Specifically, once  $\gamma_{reg}$  is set to a valid value, the BCU  
567 consistently generate a "teddy bear" with various  $\gamma_{attn}$ , otherwise it generates background or  
568 imperfect objects. This result emphasizes the role of the introduced prior energy  $E(\mathbf{K})$ . Furthermore,  
569 the  $\gamma_{attn}$  also affects to the context alignment of the generated sample (for instance  $\gamma_{attn} = 0.025$   
570 and  $\gamma_{reg} = 0.02$ ), which highlights the importance of the introduced conditional energy function  
571  $E(\mathbf{Q}; \mathbf{K})$ . The same evidences could be found in the first row in Figure 9 which are the multi-concept  
572 image generation examples.

573 **Synergy between BCU and CACAO.** While both BCU and CACAO are designed from the common  
574 energy-based perspective, each operation is originated from different energy functions  $E(\mathbf{K}; \mathbf{Q})$  and  
575  $\hat{E}(\mathbf{Q}; \{\mathbf{K}_s\}_{s=1}^M)$ , respectively. This fact suggests the synergistic energy minimization by combining  
576 the BCU and CACAO, which could further improve the text-conditional image generation. To investi-  
577 gate this further, we conducted an ablation study using a real image editing application. Specifically,  
578 we compared the editing performance when solely utilizing CACAO and when combining BCU with  
579 CACAO. The second row in Figure 9 is the result of the ablation study that shows fully-compatibility  
580 of the BCU and CACAO. Importantly, the incorporation of the BCU improves the quality of the  
581 generated images. While the CACAO alone effectively captures the context of the given editing  
582 concept, the addition of BCU enhances the fine-grained details in the generated outputs.

583 **Importance of concept negation.** Remark that a negative  $\alpha_s$  in (16) denotes the negation of given  
584 editing prompt. We empirically observed that the concept negation may significantly contribute to  
585 the performance of compositional generation. Specifically, for the image-to-image translation task in  
586 Table 1, we apply both positive and negative guidance with the target (e.g. Dog) and source (e.g. Cat)  
587 concepts, respectively, following the degree of guidance denoted in Table 2. The third row in Figure 9  
588 shows the impacts of source concept negation in the image-to-image translation task. While the  
589 positive guidance alone may fail to remove the source-concept-related features, e.g. eyes of the Cat,  
590 the negative guidance removes such conflicting existing attributes. This implies that the proposed  
591 framework enables useful arithmetic of multiple concepts for both real and synthetic image editing.

592 **Prior energy and  $\alpha$ .** While  $\frac{\alpha}{2} \text{diag}(\mathbf{K}\mathbf{K}^T)$  in (7) penalizes norm of each context vectors uniformly,  
593 the proposed prior energy function  $E(\mathbf{K})$  adaptively regularizes the smooth maximum of  $\|\mathbf{k}_i\|$ .  
594 Intuitively, adaptive penalization prevents the excessive suppression of context vectors, potentially



595 resulting in images that are more semantically aligned with a given context. To demonstrate the  
 596 effectiveness of adaptive penalization in the prior energy function, we conducted a multi-concept  
 597 image generation task with varying  $\alpha$  in (20) from 0 to 1, while fixing other hyperparameters.  
 598 Figure 10 illustrates the gradual disappearance of salient contextual elements in the generated images  
 599 depending on the change of  $\alpha$ . Specifically, the crown is the first to diminish, followed by subsequent  
 600 context elements, with the lion being the last to vanish with  $\alpha = 1$ . This result highlights the validity  
 601 of the adaptive penalization for the context vectors which stems from the prior energy function.

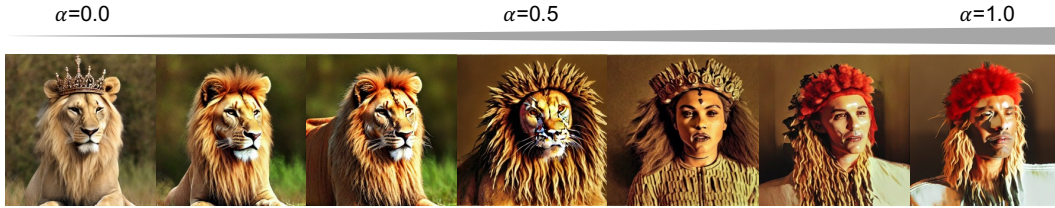


Figure 10: Generated samples with varying  $\alpha$  values. As  $\alpha$  increases, the generated images progressively deviate from the intended context, "A lion and a crown".

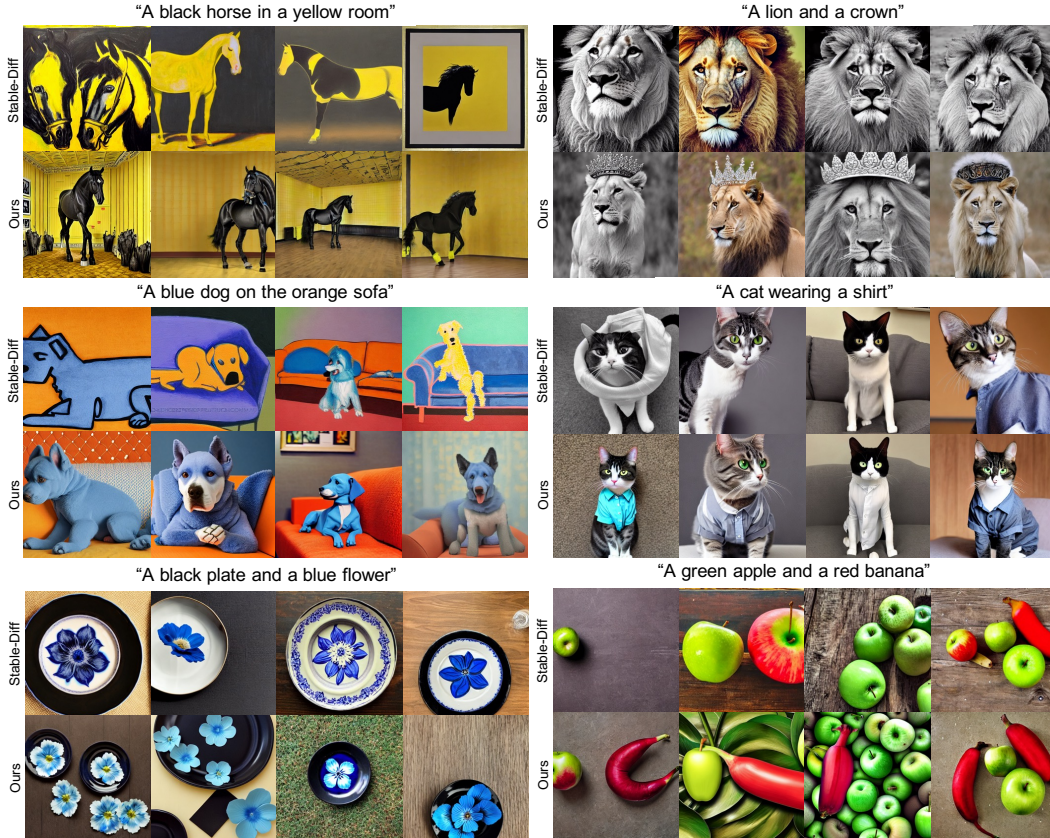


Figure 11: Further results for multi-concept image generation. Best views are displayed.



Figure 12: Further results for real image editing: horse to zebra.





Figure 13: Further results for real image editing: cat to dog.

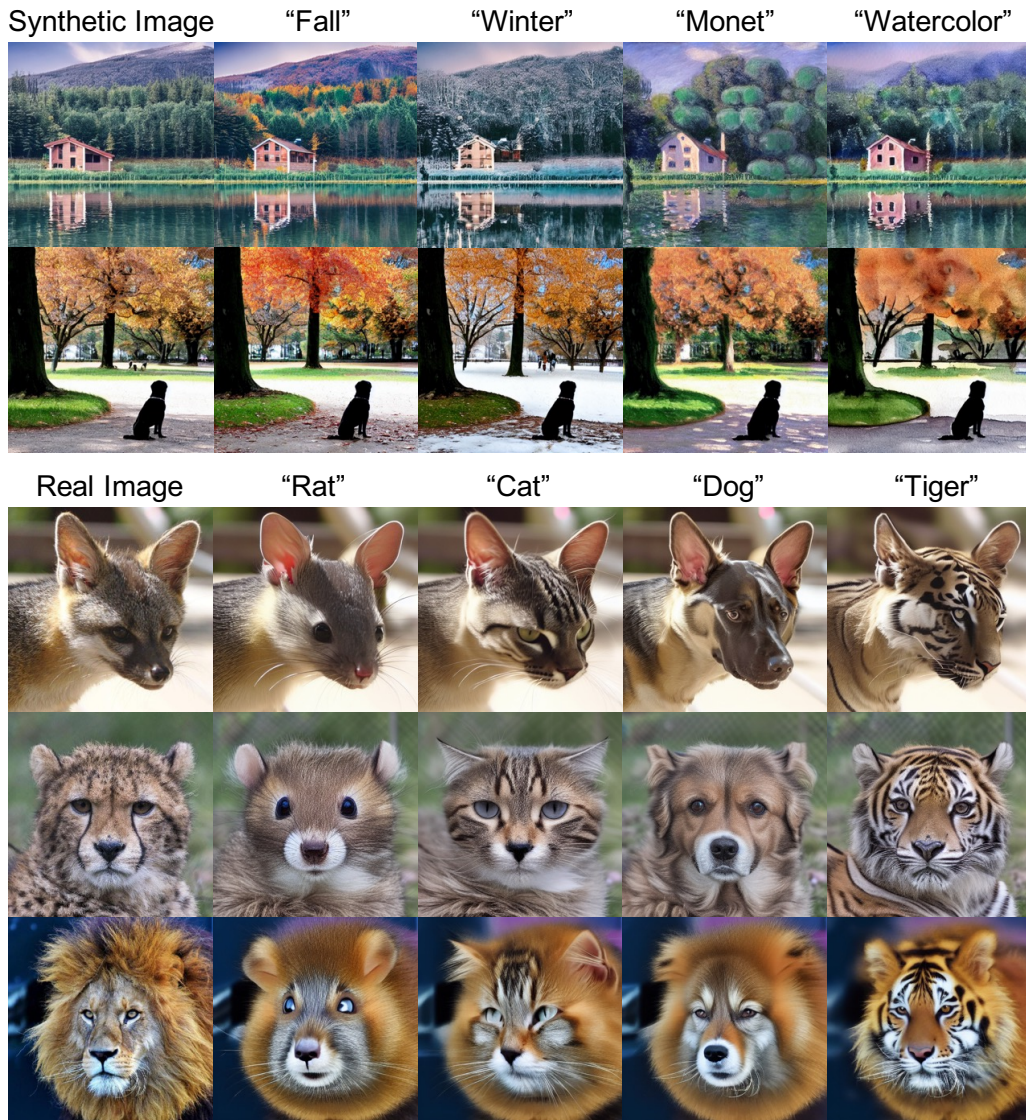


Figure 14: Further results for image editing with varying text prompts. Best views are displayed.

602 **References**

603 [1] M. Brack, F. Friedrich, D. Hintersdorf, L. Struppek, P. Schramowski, and K. Kersting. Sega: Instructing  
604 diffusion using semantic dimensions. *arXiv preprint arXiv:2301.12247*, 2023.

605 [2] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry,  
606 A. Askeel, et al. Language models are few-shot learners. *Advances in neural information processing  
607 systems*, 33:1877–1901, 2020.

608 [3] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging properties in  
609 self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer  
610 vision*, pages 9650–9660, 2021.

611 [4] H. Chefer, Y. Alaluf, Y. Vinker, L. Wolf, and D. Cohen-Or. Attend-and-excite: Attention-based semantic  
612 guidance for text-to-image diffusion models. *arXiv preprint arXiv:2301.13826*, 2023.

613 [5] Y. Choi, Y. Uh, J. Yoo, and J.-W. Ha. Stargan v2: Diverse image synthesis for multiple domains. In  
614 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8188–8197,  
615 2020.

- 616 [6] M. Demircigil, J. Heusel, M. Löwe, S. Uppgang, and F. Vermet. On a model of associative memory with  
617 huge storage capacity. *Journal of Statistical Physics*, 168:288–299, 2017.
- 618 [7] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image  
619 database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee,  
620 2009.
- 621 [8] Y. Du, S. Li, and I. Mordatch. Compositional visual generation with energy based models. *Advances in  
622 Neural Information Processing Systems*, 33:6637–6647, 2020.
- 623 [9] Y. Du, S. Li, J. Tenenbaum, and I. Mordatch. Improved contrastive divergence training of energy based  
624 models. *arXiv preprint arXiv:2012.01316*, 2020.
- 625 [10] Y. Du and I. Mordatch. Implicit generation and generalization in energy-based models. *arXiv preprint  
626 arXiv:1903.08689*, 2019.
- 627 [11] W. Feng, X. He, T.-J. Fu, V. Jampani, A. Akula, P. Narayana, S. Basu, X. E. Wang, and W. Y. Wang.  
628 Training-free structured diffusion guidance for compositional text-to-image synthesis. *arXiv preprint  
629 arXiv:2212.05032*, 2022.
- 630 [12] A. Hertz, R. Mokady, J. Tenenbaum, K. Aberman, Y. Pritch, and D. Cohen-Or. Prompt-to-prompt image  
631 editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022.
- 632 [13] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information  
633 Processing Systems*, 33:6840–6851, 2020.
- 634 [14] B. Hoover, Y. Liang, B. Pham, R. Panda, H. Strobel, D. H. Chau, M. J. Zaki, and D. Krotov. Energy  
635 transformer. *arXiv preprint arXiv:2302.07253*, 2023.
- 636 [15] T. Karras, M. Aittala, T. Aila, and S. Laine. Elucidating the design space of diffusion-based generative  
637 models. In *Proc. NeurIPS*, 2022.
- 638 [16] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks.  
639 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410,  
640 2019.
- 641 [17] D. P. Kingma, T. Salimans, B. Poole, and J. Ho. Variational diffusion models. *arXiv preprint  
642 arXiv:2107.00630*, 2021.
- 643 [18] D. Krotov and J. Hopfield. Dense associative memory is robust to adversarial inputs. *Neural computation*,  
644 30(12):3151–3167, 2018.
- 645 [19] J. Li, D. Li, C. Xiong, and S. Hoi. Blip: Bootstrapping language-image pre-training for unified vision-  
646 language understanding and generation. In *International Conference on Machine Learning*, pages 12888–  
647 12900. PMLR, 2022.
- 648 [20] L. Liu, Y. Ren, Z. Lin, and Z. Zhao. Pseudo numerical methods for diffusion models on manifolds. *arXiv  
649 preprint arXiv:2202.09778*, 2022.
- 650 [21] N. Liu, S. Li, Y. Du, A. Torralba, and J. B. Tenenbaum. Compositional visual generation with composable  
651 diffusion models. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October  
652 23–27, 2022, Proceedings, Part XVII*, pages 423–439. Springer, 2022.
- 653 [22] A. Lugmayr, M. Danelljan, A. Romero, F. Yu, R. Timofte, and L. Van Gool. Repaint: Inpainting using  
654 denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
655 and Pattern Recognition*, pages 11461–11471, 2022.
- 656 [23] C. Meng, Y. Song, J. Song, J. Wu, J.-Y. Zhu, and S. Ermon. Sdedit: Image synthesis and editing with  
657 stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021.
- 658 [24] R. Mokady, A. Hertz, K. Aberman, Y. Pritch, and D. Cohen-Or. Null-text inversion for editing real images  
659 using guided diffusion models. *arXiv preprint arXiv:2211.09794*, 2022.
- 660 [25] W. Nie, A. Vahdat, and A. Anandkumar. Controllable and compositional generation with latent-space  
661 energy-based models. *Advances in Neural Information Processing Systems*, 34:13497–13510, 2021.
- 662 [26] E. Nijkamp, M. Hill, S.-C. Zhu, and Y. N. Wu. Learning non-convergent non-persistent short-run mcmc  
663 toward energy-based model. *Advances in Neural Information Processing Systems*, 32, 2019.

- 664 [27] G. Parmar, K. K. Singh, R. Zhang, Y. Li, J. Lu, and J.-Y. Zhu. Zero-shot image-to-image translation. *arXiv preprint arXiv:2302.03027*, 2023.  
665
- 666 [28] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin,  
667 J. Clark, et al. Learning transferable visual models from natural language supervision. In *International  
668 conference on machine learning*, pages 8748–8763. PMLR, 2021.
- 669 [29] H. Ramsauer, B. Schäfl, J. Lehner, P. Seidl, M. Widrich, T. Adler, L. Gruber, M. Holzleitner, M. Pavlović,  
670 G. K. Sandve, et al. Hopfield networks is all you need. *arXiv preprint arXiv:2008.02217*, 2020.
- 671 [30] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with  
672 latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
673 Recognition*, pages 10684–10695, 2022.
- 674 [31] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation.  
675 In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International  
676 Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer,  
677 2015.
- 678 [32] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes,  
679 B. Karagol Ayan, T. Salimans, et al. Photorealistic text-to-image diffusion models with deep language  
680 understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022.
- 681 [33] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta,  
682 C. Mullis, M. Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation  
683 image-text models. *arXiv preprint arXiv:2210.08402*, 2022.
- 684 [34] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*,  
685 2020.
- 686 [35] Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in  
687 neural information processing systems*, 32, 2019.
- 688 [36] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative  
689 modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- 690 [37] N. Tumanyan, O. Bar-Tal, S. Bagon, and T. Dekel. Splicing vit features for semantic appearance transfer.  
691 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10748–  
692 10757, 2022.
- 693 [38] N. Tumanyan, M. Geyer, S. Bagon, and T. Dekel. Plug-and-play diffusion features for text-driven image-  
694 to-image translation. *arXiv preprint arXiv:2211.12572*, 2022.
- 695 [39] S. Xie, Z. Zhang, Z. Lin, T. Hinz, and K. Zhang. Smartbrush: Text and shape guided object inpainting  
696 with diffusion model. *arXiv preprint arXiv:2212.05034*, 2022.
- 697 [40] A. L. Yuille and A. Rangarajan. The concave-convex procedure (cccp). *Advances in neural information  
698 processing systems*, 14, 2001.