MITIGATING COMPOSITIONAL ISSUES IN TEXT-TO-IMAGE GENERATIVE MODELS VIA ENHANCED TEXT EMBEDDINGS

Anonymous authors

000

001

002 003 004

005

006 007 008

010 011

012

013

014

015

016

017

018

019

021

023

024

027

029

Paper under double-blind review

ABSTRACT

Text-to-image diffusion-based generative models have the stunning ability to generate photo-realistic images and achieve state-of-the-art low FID scores on challenging image generation benchmarks. However, one of the primary failure modes of these text-to-image generative models is in composing attributes, objects, and their associated relationships accurately into an image. In our paper, we investigate this compositionality-based failure mode and highlight that imperfect text conditioning with CLIP text-encoder is one of the primary reasons behind the inability of these models to generate high-fidelity compositional scenes. In particular, we show that (i) there exists an optimal text-embedding space that can generate highly coherent compositional scenes showing that the output space of the CLIP text-encoder is sub-optimal, and (ii) the final token embeddings in CLIP are erroneous as they often include attention contributions from unrelated tokens in compositional prompts. Our main finding shows that the best compositional improvements can be achieved (without harming the model's FID score) by fine-tuning *only* a simple and parameter-efficient linear projection on CLIP's representation space in Stable-Diffusion variants using a small set of compositional image-text pairs. This result demonstrates that the sub-optimality of the CLIP's output space is a major error source. We also show that re-weighting the erroneous attention contributions in CLIP can lead to slightly improved compositional performances.

1 INTRODUCTION

Text-to-image diffusion-based generative models (Rombach et al., 2021; Podell et al., 2023; Ramesh et al., 2021; Saharia et al., 2022) have achieved photo-realistic image generation capabilities on user-defined text prompts. However recent works (Huang et al., 2023) have designed compositionality benchmarks to show that these text-to-image models have low fidelity to simple compositionality prompts such as those consisting of attributes, objects, and their associated relations (e.g., "*a red book and a yellow vase*"). This hinders the use of these generative models in various creative scenarios where the end-user wants to generate a scene where the composition is derived from words (and their relationships) in the prompt.

Existing works (Chefer et al., 2023a; Feng et al., 2023; Agarwal et al., 2023; Wang et al., 2023) propose
various ways to improve compositionality in text-to-image models. These works primarily rely on modifying
the cross-attention maps by leveraging bounding box annotations and performing a small optimization in
the latent space during inference. Recent methods based on fine-tuning (Huang et al., 2023) the UNet also
lead to improved compositionality. Despite the progress, the *core reasons* behind compositionality failures in
text-to-image models remain unclear. Understanding these reasons helps designing effective methods that can
augment text-to-image models with improved compositional capabilities.

In our paper, we investigate possible reasons behind compositionality failures in text-to-image generative models. We identify two sources of errors: (i) We observe that output token embeddings in CLIP have



Figure 1: Overview of our analysis and proposed methods. The figure identifies two sources of errors in
 Stable Diffusion's inability to generate compositional prompts: (i) erroneous attention contribution in CLIP
 (minor) and (ii) sub-optimal CLIP text embedding (major). We propose a window-based linear projection
 (WiCLP), applying linear projection to a token's surrounding window to enhance embeddings.

068 significant attention contributions from irrelevant tokens, thereby introducing errors in generation. We then 069 compare the internal attention contributions in CLIP for compositional prompts to the T5 text-encoder which 070 has been shown to display strong compositional capabilities in DeepFloyd^1 . We quantitatively find that the T5 071 text-encoder displays significantly lesser erroneous attention contributions than CLIP, highlighting a potential 072 reason towards its improved compositionality. (ii) Sub-optimality of CLIP output space on compositional 073 prompts: We observe that optimizing the text embeddings, while utilizing a frozen Stable-Diffusion UNet, 074 effectively generates images with compositional scenes. We find out that there exists a text-embedding space 075 capable of generating highly coherent images with compositional scenes for various attributes (e.g., color, 076 texture, shape) which highlights that the existing CLIP output space is sub-optimal. These results indicate 077 that the output space of the CLIP text-encoder could be further improved to enable text-to-image models to generate more accurate compositional scenes. 078

079 Leveraging our observations on the deficiencies of the CLIP output space, we show that we can improve 080 the output space of the CLIP text-encoder to better align with the optimal space by applying a *simple* linear 081 projection on top of CLIP (see Figure 1). This leads to stronger compositional performances. In particular, 082 we propose Window-based Compositional Linear Projection (WiCLP), a lightweight fine-tuning method that 083 significantly improves the model's performance on compositional prompts, yielding results comparable to existing baselines (see Figure 2). Moreover, it preserves the model's clean accuracy, as evidenced by a low 084 FID on clean prompts, offering a *parameter* and *speed*-efficient solution. We also show that reweighting the erroneous attention contributions in CLIP can lead to improved compositional performances, however, the improvements often lag behind WiCLP. 087

Fine-tuning a subset of components of the diffusion model can result in an increase in the FID score for clean prompts. While fine-tuning only a linear projection partially mitigates this, we find that applying it over all the time steps results in an increase in FID. To mitigate this, we introduce SWITCH-OFF where we only apply
 WiCLP during the initial steps of generation, switching it off for the remaining steps. This enables the model

092

¹https://huggingface.co/DeepFloyd/IF-I-M-v1.0

094		SD v1.4	CLP	SD v2	WiCLP	
095						
096						
097	A blue backpack and					
098	a red chair				- Andrew -	
099						
100						
101			00	1000		
102		V.C.P.R.				
103	A yellow book and					
104	a leu vase					
105						
106		Nykolast 1/1-				
107	Figure 2: Quali	itativa comparison k	atwaan CID a	nd WiCID ve th	a basalinas	
108	Figure 2. Quan	native comparison t			e basennes.	
109						
110	to obtain a coherent compositi	onal scene in early	steps (crucial t	for compositiona	al prompts) while	retaining
111	clean accuracy on surrounding	prompts, as the gen	eration in final	steps is guided b	by the original tex	t-encoder
112	not the augmented one that ma	ps to the optimized	space.	r c		
113	In summary our contributions	are as follows:				
114	in summary, our controlations	die ds follows.				
115	• We perform an in-de	pth analysis of the	reasons behind	d compositional	ity failures in ope	en-source
110	text-to-image generat	ive models, highligl	hting two reaso	ons for them.		
118	• Leveraging our obser	rvations, we propos	se WiCLP for	Stable Diffusior	v-1.4 and v-2 w	hich can
110	augment the models	with improved con	npositionality	while preservin	g their clean acc	uracy on
120	surrounding prompts	. We observe impro	ovements of 16	5.18%, 15.15%, a	and 9.51% on SD	v1.4 and
121	14.35%, 11.14%, and	6% on SD v2 in V	QA scores (Hu	uang et al., 2023) across color, tex	ture, and
122	shape datasets, respe	ctively. Our metho	od achieves co	mpetitive VQA	scores compared	l to other
123	baselines, while dem	ionstrating superior	FID on clear	n prompts, requi	ring fewer param	<i>ieters</i> for
124	opunization, and ena	bing jast injerence				
125	Overall, our paper provides qu	antitative evidence	e elucidating th	ne compositional	challenges with	n text-to-
126	image models and strong basel	lines to mitigate suc	h issues.	I	0	
127	c c	C				
128	$2 B_{ACKCPOUND}$					
129	2 DACKGROUND					
130	Compositionality in Tayt to 1	maga Concrativa	Modola Am	agent month Illuon	a = at al (2022); est	na du a a a
131	benchmark for testing composit	tionality in text-to-ir	nage models sh	owing the suscer	g et al. (2025) Illu	urce text-
132	to-image models on simple com	positional prompts	In addition th	e authors also pro	prose a fine-tuning	baseline
133	to augment text-to-image mod	lels with improved	compositional	ity. The composition	itionality issue ca	n also be
134	addressed at inference time by	modifying the cross	-attention map	s leveraging hand	1-crafted loss func	tions and
135	bounding boxes generated fro	m a language mode	el (Chefer et al	l., 2023a; Feng e	t al., 2023; Agary	wal et al.,
136	2023; Wang et al., 2023; Nie e	t al., 2024; Lian et a	al., 2023; Liu e	et al., 2022a). Ho	wever, Huang et a	al. (2023)
137	show that a data-driven and fin	e-tuning approach i	is more suitabl	e towards improv	ving compositiona	ality.
138						
139	Interpretability of Text-to-In	nage Generative N	Models. The	re have been reco	ent efforts to inter	pret text-
140	to-image models like Stable D	ittusion. DAAM (T	ang et al., 202	3; Hertz et al., 20	J22) studies the g	eneration

process in diffusion models by analyzing cross-attention maps between text tokens and image pixels, high lighting their semantic precision. Basu et al. (2023) use causal tracing to understand how knowledge is stored
 in models like Stable Diffusion v1 while Rezaei et al. (2024) propose a mechanistic approach to localize
 knowledge in cross-attention layers of various text-to-image models. Chefer et al. (2023b) explore concept
 decomposition in diffusion models.

2.1 TEXT-TO-IMAGE DIFFUSION MODELS: TRAINING AND INFERENCE

In diffusion models, noise is added to the data following a Markov chain across multiple time-steps $t \in [0, T]$. 149 Starting from an initial random real image \mathbf{x}_0 along with its caption c, $(\mathbf{x}_0, c) \sim \mathcal{D}$, the noisy image at 150 time-step t is defined as $\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{(1-\alpha_t)} \epsilon$. The denoising network denoted by $\epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}, t)$ is 151 pre-trained to denoise the noisy image x_t to obtain x_{t-1} . For better training efficiency, the noising along 152 with the denoising operation occurs in a latent space defined by $\mathbf{z} = \mathcal{E}(\mathbf{x})$, where \mathcal{E} is an encoder such as 153 VQ-VAE (van den Oord et al., 2017). Usually, the conditional input c to the denoising network $\epsilon_{\theta}(.)$ is a 154 text-embedding of the caption c through a text-encoder $\mathbf{c} = v_{\gamma}(c)$. The pre-training objective for diffusion 155 models can be defined as follows: 156

$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{x}_0, c) \sim \mathcal{D}, \epsilon, t} \left[\left\| \epsilon - \epsilon_{\theta}(\mathbf{z}_t, \mathbf{c}, t) \right\|_2^2 \right],$$

where θ is the set of learnable parameters in the UNet ϵ_{θ} . During inference, where the objective is to synthesize an image given a text-embedding c, a random Gaussian noise $\mathbf{z}_T \sim \mathcal{N}(0, I)$ is iteratively denoised for a fixed range of time-steps to produce the final image.

2.2 COMPOSITIONALITY EVALUATION METRICS

We focus on the disentangled BLIP-Visual Question Answering (referred to as VQA for simplicity) score proposed by Huang et al. (2023) as a key metric for evaluating image quality. The VQA score measures how accurately an image captures the compositional elements described in the prompt, offering a closer correlation with human judgment compared to metrics like CLIP-Score (Hessel et al., 2021).

169 170 2.3 DATASET COLLECTION

171 We utilize the T2I-CompBench dataset (Huang et al., 2023), focusing on three key categories: color, texture, 172 and shape, with a total of 1,000 prompts across both training and evaluation sets. T2I-CompBench is a well-173 established and widely recognized dataset (Esser et al., 2024). This dataset provides distinct training and 174 evaluation splits for each category, enabling a structured approach to assessing performance. To generate 175 high-quality images, we use three generative models: SD 1.4 (Rombach et al., 2021), DeepFloyd, and SynGen 176 (Rassin et al., 2024), creating 100 samples per prompt with SD 1.4, 60 with DeepFloyd, and 50 with SynGen. This ensures a wide variety of generated images, leveraging each model's strengths. For each prompt, we 177 combined all 210 samples from the three models and selected the top 30 with the highest VQA scores, 178 ensuring the final dataset consisted of images that most accurately reflected the prompts. 179

180 181

182

147

148

157 158 159

160

161 162 163

164

3 SOURCE (I) : ERRONEOUS ATTENTION CONTRIBUTIONS IN CLIP

In this section, we leverage attention contributions (Elhage et al., 2021; Dar et al., 2023) to analyze the text-embeddings of compositional prompts in the CLIP text-encoder (which is commonly used in many text-to-image models) and compare them with T5-text encoder of DeepFloyd, a model which results in stronger compositionality. Many of the compositional prompts from Huang et al. (2023) have a decomposable template of the form $\mathbf{a}_i \mathbf{o}_j + \mathbf{a}_j \mathbf{o}_j$, where $\mathbf{a}_i, \mathbf{a}_j$ are attributes (e.g., "black", "matted") while $\mathbf{o}_i, \mathbf{o}_j$ describe



203

204

205

206 207

218

220 221 222



the more accurate performance of T5.

Figure 3: The heatmap illustrates unintended at- Figure 4: Quantitatively, we find CLIP to have signifitention contributions in CLIP, while highlighting cantly higher erroneous attention contributions averaged across 780 prompts of color dataset and 582 prompts of texture dataset.

208 objects (e.g., "car", "bag"). We use attention contributions to understand how the text-embeddings of the 209 compositional tokens (e.g., $\mathbf{a}_i, \mathbf{a}_i, \mathbf{o}_i, \mathbf{o}_i$) are formed for both T5 and CLIP over the layers of these models.

210 The attention mechanism in layer ℓ of a transformer consists of four weight matrices 211 W_q, W_v, W_k, W_o (Vaswani et al., 2017). Each of these weight matrices is divided into H heads de-212 noted by $W_a^h, W_v^h, W_k^h \in \mathbb{R}^{d \times d_h}, W_o^h \in \mathbb{R}^{d_h \times d}$ for all $h \in [H]$. Note that d_h is the dimension of the 213 internal token embeddings. We omit ℓ for simplicity, but each layer has its own attention matrices. These 214 matrices are applied on the token embeddings of the output of layer $\ell - 1$, denoted by $\bar{\mathbf{x}}_i$ for token j in that 215 layer. We denote by q_i^h, k_i^h , and v_i^h the projection of $\bar{\mathbf{x}}_i$ on query, key, and value matrices of the h-th head of 216 layer ℓ . More precisely, 217

$$\mathbf{q}_j^h = \bar{\mathbf{x}}_j W_{\mathbf{q}}^h, \quad \mathbf{k}_j^h = \bar{\mathbf{x}}_j W_{\mathbf{k}}^h, \quad \mathbf{v}_j^h = \bar{\mathbf{x}}_j W_{\mathbf{v}}^h.$$

The *contribution* of token j to token i in layer ℓ , denoted by cont_{i,j}, is computed as follows: 219

$$\operatorname{cont}_{i,j} = \left\| \sum_{h=1}^{H} \operatorname{attn}_{i,j}^{h} \operatorname{v}_{j}^{h} W_{o}^{h} \right\|$$

where attn^h_{i,j} is the attention weight of token i to j in the h-th head of layer ℓ . Specifically,

$$\operatorname{attn}_{i,.}^{h} = \operatorname{SOFTMAX}\left(\left\{\frac{\langle \mathbf{q}_{i}^{h}, \mathbf{k}_{j}^{h} \rangle}{\sqrt{d_{h}}}\right\}_{j=1}^{n}\right).$$

227 Notably, $cont_{i,j}$ is a significant metric that quantifies the *contribution* of a token j to the norm of a token i 228 at layer ℓ . We employ this metric to identify layers in which important tokens highly attend to *unintended* 229 tokens, or lowly attend to *intended* ones. We refer to Appendix C.1 for more details on attention contribution. 230

231 3.1 KEY FINDING: T5 HAS LESS ERRONEOUS ATTENTION CONTRIBUTIONS THAN CLIP 232

233 We refer to Figure 3 that visualizes attention contribution of both T5 and CLIP text-encoder in the last layer $(\ell = 11)$ for the prompt "a green bench and a red car". Ideally, the attention mechanism should guide the token

"car" to focus more on "red" than "green", but in the last layer of the CLIP text-encoder, "car" significantly attends to "green". In contrast, T5 shows a more consistent attention pattern, with "red" contributing more to the token "car" and "green" contributing more to the token "bench".

We further conduct an extensive analysis on specific types of prompts, consisting of 780 prompts of color 239 dataset and 582 prompts of texture dataset, each structured as " $\mathbf{a}_1 \mathbf{o}_1$ and $\mathbf{a}_2 \mathbf{o}_2$." For each prompt, we obtain 240 attention contributions in all layers and count the number of layers where unintended attention contributions 241 occur. In the CLIP text-encoder, unintended attention occurs when o_2 attends more to a_1 than a_2 . For T5, it 242 occurs when o_2 attends more to a_1 than a_2 , or o_1 attends more to a_2 than a_1 . Figure 4 provides a quantitative 243 comparison of unintended attention across various prompts between the CLIP text-encoder and T5. The T5 244 model demonstrates improved performance on our metric compared to the CLIP text-encoder, reinforcing the 245 hypothesis that erroneous attention mechanisms in CLIP may contribute to its weaker compositionality in 246 text-to-image models. This aligns with the general observation that pretrained text-to-image models using the T5 text-encoder tend to exhibit superior compositionality. Additional details can be found in Appendix C.4. 247 Further experiments with other text-encoders are also reported in Appendix C.3. 248

250 3.2 ZERO-SHOT ATTENTION REWEIGHTING

Inspired by attention mechanism shortcomings of CLIP text-encoder, we aim to improve compositionality of CLIP-based diffusion models by zero-shot reweighting of the attention maps. Specifically, we apply a hand-crafted zero-shot manipulation of the attention maps in certain layers of the CLIP text-encoder to effectively reduce unintended attentions while enhancing meaningful ones. This zero-shot reweighting is applied to the logits before the SOFTMAX layer in the last three layers of the text-encoder. More precisely, we compute a matrix $M \in \mathbb{R}^{n \times n}$ and add it to the attention logits. For each head h, the new attention values are computed and then propagated through the subsequent layers of the text encoder:

$$\operatorname{attn}_{i,.}^{'h} = \operatorname{SOFTMAX}\left(\left\{\frac{\langle \mathbf{q}_{i}^{h}, \mathbf{k}_{j}^{h}\rangle}{\sqrt{d_{h}}} + M_{i,j}\right\}_{j=1}^{n}\right)$$

We set the values in M by considering the ideal case where no incorrect attentions occur in the mechanism. For example, for prompt "a green bench and a red car", we ensure that the token "car" does not attend to the token "green" by assigning a sufficiently large negative value to the corresponding entry in matrix M. Further details on how we obtain matrix M can be found in Appendix C.2.

Key Results. Applying zero-shot attention reweighting with matrix M on 780 compositional prompts of color dataset, we achieved a 2.93% improvement in VQA scores. Examples of effective zero-shot reweighting, demonstrating its impact on mitigating compositionality issues in can be found in Appendix C.2. Although erroneous attention contributions in the CLIP text-encoder is one source of error, it is not the primary error source due to modest improvements in compositional accuracy. In the next section, we investigate the sub-optimality of the output space of CLIP text-encoder, which we find to be a significant source of error.

4 SOURCE (II) : SUB-OPTIMALITY OF CLIP TEXT-ENCODER FOR COMPOSITIONAL PROMPTS

In this section, we understand if the UNet is capable of generating compositional scenes by optimizing the textembeddings that it takes as the conditional input. Given an input prompt c with a particular composition (e.g., "*a red book and a yellow table*"), we utilize our dataset and obtain \mathcal{D}_c including high-quality compositional images for prompt c. We then optimize the output text-embedding \mathbf{c} as follows:

$$\mathbf{c}^* = rg\min_{\mathbf{c}} \mathbb{E}_{x_0 \sim \mathcal{D}_c, \epsilon, t} \left[\|\epsilon - \epsilon_{ heta}(\mathbf{z}_t, \mathbf{c}, t)\|_2^2
ight].$$

280 281

272 273

274

275

249

Figure 5: Comparative analysis of VQA Scores between CLIP text-embeddings and optimized textembeddings using Stable Diffusion v1.4 across color, texture, and shape categories. Results show CLIP text embeddings achieved scores of 0.3615 for color, 0.4306 for texture, and 0.3619 for shape, while optimized text embeddings achieved scores of 0.7513 for color, 0.7254 for texture, and 0.58728 for shape.

We then use c^* to generate images using the UNet ϵ_{θ} across different seeds. Figure 5 depicts a few of generated images using optimized text-embeddings.

Key Results. As seen in Figure 5, we consistently improve VQA scores across a variety of compositional prompts (i.e., color, texture, and shape). This indicates that CLIP text-encoder does not output the proper text-embedding suitable for generating compositional scenes. However, that optimized embedding space exists, highlighting the ability of UNet to generate coherent compositional scenes when a proper text-embedding is given. This further motivates the idea of improving CLIP output space to mitigate compositionality issues in text-to-image diffusion models. We refer to Appendix B for other configurations showing that optimizing a subset of tokens can also effectively improve compositionality.

5 LINEAR PROJECTION ON CLIP: A SIMPLE BASELINE TO IMPROVE COMPOSITIONALITY IN TEXT-TO-IMAGE GENERATIVE MODELS

In this Section, we provide two baselines CLP and WiCLP that are linear modification of CLIP output to map that sub-optimal space to an enhanced one, better suited for compositionality.

5.1 CLP: TOKEN-WISE COMPOSITIONAL LINEAR PROJECTION

Given the text-embedding $\mathbf{c} \in \mathbb{R}^{n \times d}$ as the output of the text-encoder for prompt c, i.e., $\mathbf{c} = v_{\gamma}(c)$, we train a linear projection $\operatorname{CLP}_{W,b} : \mathbb{R}^{n \times d} \to \mathbb{R}^{n \times d}$. This projection includes a matrix $W \in \mathbb{R}^{d \times d}$ and a bias term $b \in \mathbb{R}^d$, which are applied token-wise to the output text-embeddings of the encoder. More formally, for $\mathbf{c} \in \mathbb{R}^{n \times d}$ including text-embeddings of n tokens $\mathbf{c}_1, \mathbf{c}_2, \cdots, \mathbf{c}_n \in \mathbb{R}^d$, $\operatorname{CLP}_{W,b}(\mathbf{c})$ is obtained by stacking projected embeddings $\mathbf{c}'_1, \mathbf{c}'_2, \cdots, \mathbf{c}'_n$ where $\mathbf{c}'_i = W^T \mathbf{c}_i + b$.

Finally, we solve the following optimization problem on a dataset \mathcal{D} including image-caption pairs of high-quality compositional images:

$$W^{*}, b^{*} = \arg\min_{W, b} \mathbb{E}_{(x_{0}, c) \sim \mathcal{D}, \epsilon, t} \left[\left\| \epsilon - \epsilon_{\theta} \left(\mathbf{z}_{t}, \mathtt{CLP}_{W, b} \left(\mathbf{c} \right), t \right) \right\|_{2}^{2} \right]$$

328 We then apply CLP_{W^*,b^*} on CLIP text-encoder to obtain improved embeddings.



282

286

287

288

289 290

291

292

293 294

295

296

297

298 299

309 310

311

312 313

314

315 316

317





Figure 7: Trade-off between VQA and FID scores with SWITCH-OFF at different thresholds.

5.2 WICLP: WINDOW-BASED COMPOSITIONAL LINEAR PROJECTION

353 354 355

356

365 366

369 370

In this section, we propose a more advanced linear projection scheme where the new embedding of a token
 is derived by applying a linear projection on that token in conjunction with a set of its adjacent tokens, i.e.,
 tokens within a specified window. This method not only leverages the benefits of CLP but also incorporates
 the contextual information from neighboring tokens, potentially leading to more precise text-embeddings.

More formally, we train a mapping $\operatorname{WiCLP}_{W,b} : \mathbb{R}^{n \times d} \to \mathbb{R}^{n \times d}$ including a parameter *s* (indicating window length), matrix $W \in \mathbb{R}^{(2s+1)d \times d}$, and a bias term $b \in \mathbb{R}^d$. For text-embeddings $\mathbf{c} \in \mathbb{R}^{n \times d}$ consisting of *n* token embeddings of $\mathbf{c}_1, \mathbf{c}_2, \cdots, \mathbf{c}_n \in \mathbb{R}^d$, we obtain $\operatorname{WiCLP}_{W,b}$ by stacking projected embeddings $\mathbf{c}'_1, \mathbf{c}'_2, \cdots, \mathbf{c}'_n$ where

$$\mathbf{c}'_{i} = W^{T}$$
 Concatenation $\left((\mathbf{c}_{j})_{j=i-s}^{i+s} \right) + l$

367 Similarly, we solve the following optimization problem to train the projection: 368

$$W^{*}, b^{*} = \arg\min_{W, b} \mathbb{E}_{(x_{0}, c) \sim \mathcal{D}, \epsilon, t} \left[\left\| \epsilon - \epsilon_{\theta} \left(\mathbf{z}_{t}, \texttt{WiCLP}_{W, b} \left(\mathbf{c} \right), t \right) \right\|_{2}^{2} \right]$$

371 Note that we use s = 2, i.e., window length of 5 in our experiments.

Comparison between CLP and WiCLP. We observe that WiCLP improves over CLP (special case of WiCLP with s = 0) by incorporating adjacent tokens in addition to the actual token. This approach likely improves embeddings by mitigating unintended attention from adjacent tokens. For discussion on choosing the window length (s) in WiCLP, see Appendix D.6.

376		Color	Texture	Shape	
377		Baseline	0.3765	0.4156	0.3576
370	Stable Diffusion v1.4	CLP	0.4837	0.5312	0.4307
379		WiCLP	0.5383	0.5671	0.4527
201		Baseline	0.5065	0.4922	0.4221
202		Composable (Liu et al., 2022b)	0.4063	0.3645	0.3299
202		Structured (Feng et al., 2022)	0.4990	0.4900	0.4218
207	Stable Diffusion v2	Attn-Exct (Chefer et al., 2023a)	0.6400	0.5963	0.4517
295		GORS-unbaised (Huang et al., 2023)	0.6414	0.6025	0.4546
305		CLP	0.6075	$\bar{0.5707}$	0.4567
387		WiCLP	0.6500	0.6036	0.4821

Table 1: Quantitative comparison with state-of-the-art and baseline methods across different categories of the T2I-CompBench dataset

5.3 SWITCH-OFF: TRADE-OFF BETWEEN COMPOSITIONALITY AND CLEAN ACCURACY

Fine-tuning models or adding modules to a base model often results in a degradation of image quality and an increase in the Fréchet Inception Distance (FID) score. To balance the trade-off between improved compositionality and the quality of generated images for clean prompts – an important issue in existing work – inspired by Hertz et al. (2022), we adopt SWITCH-OFF, where we apply the linear projection only during the initial steps of inference. Specifically, given a time-step threshold τ , for $t \ge \tau$, we use WiCLP_{W*,b*}(c), while for $t < \tau$, we use the unchanged embedding c as the input to the cross-attention layers.

Figure 7 illustrates the trade-off between VQA score and FID on a randomly sampled subset of MS-COCO (Lin et al., 2014) for different choices of τ . As shown, even a large value of τ suffices for obtaining high-quality compositional scenes as the composition of final generated image is primarily formed at early steps. Thus, choosing a large τ preserves the model's improved compositionality while maintaining its clean accuracy. Setting $\tau = 800$ offers a competitive VQA score compared to the model where projection is applied at all time steps, and achieves a competitive FID similar to that of the clean model. Figure 6 depicts a few images generated using different choices of τ . We refer to Appendix D.5 for more visualizations.

408 409

388

389

390

392 393

394

6 EXPERIMENTS

410 411

Existing Baselines. We evaluate the performance of four methods alongside standard models SD v1.4 and SD v2. These include Composable Diffusion (Liu et al., 2022b), which addresses concept conjunction and negation in pretrained diffusion models; Structured Diffusion (Feng et al., 2022), which focuses on attribute binding; Attn-Exct (Chefer et al., 2023a), which ensures correct attention to all subjects in the prompt; and GORS (Huang et al., 2023), which fine-tunes Stable Diffusion v2 using a reward function. GORS optimizes more parameters but underperforms slightly compared to our method, while Attn-Exct requires iterative optimizations during inference, making it slower than our method, which adds only a linear projection layer.

Training Setup. All of the models are trained using the objective function of diffusion models on color,
 texture, and shape datasets. During training, we keep all major components frozen, including the U-Net,
 CLIP text-encoder, and VAE encoder and decoder, and only the linear projections are trained. We refer to
 Appendix D.1 for details on the training procedure.



prompts compared to base models, but this increase is smaller than other baselines—for example, WiCLP scores 27.40 versus GORS at 30.54. Further FID performance details are available in Appendix D.2.

Human Experiments. We conducted a human evaluation where participants compared images generated by
SD v1.4 and SD v1.4 + WiCLP, selecting the image that best matched the given prompt. The results showed
that in 34.625% of cases, evaluators chose the base model's image; in 51.875%, they preferred the WiCLP
images; and in 13.50%, they rated both equally. Further details can be found in Appendix D.2.

7 CONCLUSION

459 460

461

462 Our paper examines potential error sources in text-to-image models for generating images from compositional 463 prompts. We identify two error sources: (i) A minor error source, where the token embeddings in the 464 CLIP text-encoder have erroneous attention contributions and (ii) A major error source, where we find the 465 output space of the CLIP text-encoder to be sub-optimally aligned to the UNet for compositional prompts. 466 Leveraging our observations, we propose a simple and strong baseline WiCLP which involves fine-tuning a 467 linear projection on CLIP's representation space. WiCLP though inherently simple and parameter efficient, 468 outperforms existing methods on compositional image generation benchmarks and maintains a low FID score on a broader range of clean prompts. We discuss limitations in Appendix A. 469

470 REFERENCES

476

- Aishwarya Agarwal, Srikrishna Karanam, K J Joseph, Apoorv Saxena, Koustava Goswami, and Balaji Vasan
 Srinivasan. A-star: Test-time attention segregation and retention for text-to-image synthesis, 2023.
- Samyadeep Basu, Nanxuan Zhao, Vlad Morariu, Soheil Feizi, and Varun Manjunatha. Localizing and editing knowledge in text-to-image generative models, 2023.
- Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models, 2023a.
- Hila Chefer, Oran Lang, Mor Geva, Volodymyr Polosukhin, Assaf Shocher, Michal Irani, Inbar Mosseri, and
 Lior Wolf. The hidden language of diffusion models. *arXiv preprint arXiv:2306.00966*, 2023b.
- 481 Guy Dar, Mor Geva, Ankit Gupta, and Jonathan Berant. Analyzing transformers in embedding space, 2023.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda
 Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac
 Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario
 Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. https://transformercircuits.pub/2021/framework/index.html.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi,
 Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey,
 Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution
 image synthesis, 2024. URL https://arxiv.org/abs/2403.03206.
- Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Akula, Pradyumna Narayana, Sugato Basu,
 Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis. *arXiv preprint arXiv:2212.05032*, 2022.
- Weixi Feng, Wanrong Zhu, Tsu jui Fu, Varun Jampani, Arjun Akula, Xuehai He, Sugato Basu, Xin Eric
 Wang, and William Yang Wang. Layoutgpt: Compositional visual planning and generation with large
 language models, 2023.
- Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to prompt image editing with cross attention control, 2022.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith.
 Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering, 2023. URL
 https://arxiv.org/abs/2303.11897.
- Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation, 2023.
- Long Lian, Boyi Li, Adam Yala, and Trevor Darrell. Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models. *ArXiv*, abs/2305.13655, 2023. URL https://api.semanticscholar.org/CorpusID:258841035.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro
 Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in
 context. *CoRR*, abs/1405.0312, 2014. URL http://arxiv.org/abs/1405.0312.

- 517
 518
 518
 519
 519
 520
 520
 519
 520
 520
 521
 521
 522
 522
 523
 524
 524
 524
 524
 524
 525
 526
 527
 528
 529
 529
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 521
 521
 522
 522
 523
 524
 524
 524
 524
 525
 526
 527
 528
 529
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
 520
- Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *European Conference on Computer Vision*, pp. 423–439. Springer, 2022b.
- Weili Nie, Sifei Liu, Morteza Mardani, Chao Liu, Benjamin Eckart, and Arash Vahdat. Compositional
 text-to-image generation with dense blob representations, 2024.
- Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis, 2023.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya
 Sutskever. Zero-shot text-to-image generation. In Marina Meila and Tong Zhang (eds.), Proceedings of the
 38th International Conference on Machine Learning, volume 139 of Proceedings of Machine Learning
 Research, pp. 8821–8831. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/
 v139/ramesh21a.html.
- Royi Rassin, Eran Hirsch, Daniel Glickman, Shauli Ravfogel, Yoav Goldberg, and Gal Chechik. Linguistic
 binding in diffusion models: Enhancing attribute correspondence through attention map alignment, 2024.
- Keivan Rezaei, Samyadeep Basu, Ryan Rossi, Cherry Zhao, Vlad Morariu, Varun Manjunatha, and Soheil Feizi. On mechanistic knowledge localization in text-to-image generative models. *arXiv preprint arXiv:2405.01008*, 2024.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution
 image synthesis with latent diffusion models. *CoRR*, abs/2112.10752, 2021. URL https://arxiv.org/abs/2112.10752.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems, volume 35, pp. 36479–36494. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/ec795aeadae0b7d230fa35cbaf04c041-Paper-Conference.pdf.
- Raphael Tang, Linqing Liu, Akshat Pandey, Zhiying Jiang, Gefei Yang, Karun Kumar, Pontus Stenetorp, Jimmy Lin, and Ferhan Ture. What the DAAM: Interpreting stable diffusion using cross attention. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5644–5659, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.310. URL https://aclanthology.org/2023.acl-long.310.
- Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. *CoRR*, abs/1711.00937, 2017. URL http://arxiv.org/abs/1711.00937.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017. URL http://arxiv.org/abs/1706.03762.
- Ruichen Wang, Zekang Chen, Chen Chen, Jian Ma, Haonan Lu, and Xiaodong Lin. Compositional text-to image synthesis with attention map control of diffusion models, 2023.



Figure 9: Comparison of VQA scores when optimizing different subsets of tokens for the sample prompt: "A red book and a yellow vase"

A LIMITATIONS

In this paper, we have thoroughly analyzed one of the key reasons why Stable Diffusion struggles to generate compositional prompts and proposed a lightweight method to mitigate this issue. However, there remains significant room for improvement in this area. Our approach focuses on improving the text encoder, which we identified as a major source of error. There are potentially other sources of the issue within the entire generative model pipeline that need to be explored. Additionally, our method involves a small fine-tuning step using a simple linear projection. Future work could explore alternative approaches, such as more sophisticated fine-tuning techniques, advanced attention mechanisms, or hybrid models that integrate multiple strategies.

B Optimizing the Text-embeddings of a Subset of Tokens

600 Given $\mathbf{c} \in \mathbb{R}^{n \times d}$, where *n* refers to the number of tokens and *d* refers to the dimensionality of the text-601 embedding, for the second configuration we only optimize a subset of tokens $n' \in n$. We refer to this subset 602 of tokens as \mathbf{c}' . These tokens correspond to relevant parts of the prompt which govern compositionality (e.g., 603 "red book" and "yellow table" in "A red book and an yellow table").

 $\mathbf{c}'^* = \arg\min_{\mathbf{c}'} \mathbb{E}_{\epsilon,t} ||\epsilon - \epsilon_{\theta}(\mathbf{z}_t, \mathbf{c}', t)||_2^2,$

Figure 9 shows the results for the sample prompt "a red book and a yellow vase". We considered different subsets of tokens n': adjectives ("red" and "yellow"), nouns ("book" and "vase"), both nouns and adjectives, and all tokens in the sentence. The results indicate that optimizing even a few tokens significantly improves the VQA score. However, optimizing all tokens in the sentence yields the highest score.



Figure 10: Visualization of attention map and attention contribution for prompt "a green bench and a red car" over different layers of CLIP. Contribution provides better insight on the attention mechanism.

C SOURCE (I) : ERRONEOUS ATTENTION CONTRIBUTIONS

C.1 ATTENTION CONTRIBUTION

In this Section, we provide more details on our analysis to quantitatively measure tokens' contribution to each other in a layer of attention mechanism. One natural way of doing this analysis is to utilize attention maps $\operatorname{attn}_{i,j}^h$ and aggregate them over heads, however, we observe that this map couldn't effectively show the contribution. Attention map does not consider norm of tokens in the previous layer, thus, does not provide informative knowledge on how each token is formed in the attention mechanism. In fact, as seen in Figure 10, we cannot obtain much information by looking at these maps while attention contribution clearly shows amount of norm that comes from each of the attended tokens.

642

624

629 630

631 632

633

634

635

637

C.2 ZERO-SHOT ATTENTION REWEIGHTING

To fix unintended attentions, we aim to compute a matrix M to be applied across various heads in the last few layers of CLIP, reducing the effect of wrong attention, leading to more accurate text-embeddings that are capable of generating high-quality compositional scenes. To avoid unintended attention for prompts of the form " $a_1o_1 + a_2o_2$ ", we add large negative values to entries M_{o_2,a_1} , M_{a_2,a_1} , and some positive value to M_{o_2,a_2} and M_{o_1,a_1} , and small negative value to M_{o_2,o_1} . To find what values to assign to those entries, we consider a small set of prompts in color dataset (5 prompts in total) and obtain parameters for that matrix to maximize VQA score. Figure 11 shows few examples of zero-shot modification.

- 650
- 651 C.3 EXPERIMENTS WITH LLAMA3 8B 652

We explored the analysis of attention contributions to identify unintended attention in LLaMa3 8B, which utilizes a more advanced text encoder specifically designed for language modeling and pretrained on largescale text corpora. Table 2 reports the rate of unintended attention across prompts in the color and texture datasets. The results demonstrate that unintended attention occurs less frequently in more advanced text encoders, further emphasizing the limitations of the CLIP text encoder.

	color		texture	
	last layer	all layers	last layer	all layers
LLaMa3	0.015	0.081	0.033	0.066
CLIP	0.657	0.187	0.696	0.213

Table 2: Unintended attention rate in LLaMa3 8B vs CLIP. LLaMa3 shows significant less unintended attentions.

	SD v1.4	SD v2	SD v1.4 + WiCLP	$SD \; v2 \text{+} \text{WiCLP}$	GORS
FID Score	24.33	23.27	25.40	27.40	30.54

Table 3: Comparison of FID scores between the baseline models and WiCLP using SWITCH-OFF with $\tau = 800$, as well as the GORS approach.

C.4 MODELS WITH T5 TEXT-ENCODER

We conducted experiments to measure the VQA score on the color dataset for models that use T5 as their text encoder. DeepFloyd achieved a score of 0.604, which is significantly higher than that of SD-v1.4.
Additionally, DeepFloyd-I-M, which employs a smaller first-stage UNet compared to DeepFloyd, obtained a score of 0.436, also surpassing the SD-v1.4 score.

D EXPERIMENTS

664

665

671 672 673

674 675

676

677

678

679 680

681 682

683

693

694

D.1 TRAINING SETUP

684 In this section, we present the details of the experiments conducted to evaluate our proposed methods. The 685 training is performed for 25,000 steps with a batch size of 4. An RTX A5000 GPU is used for training models 686 based on Stable Diffusion 1.4, while an RTX A6000 GPU is used for models based on Stable Diffusion 2. 687 We employed the Adam optimizer with a learning rate of 1×10^{-5} and utilized a Multi-Step learning rate scheduler with decays ($\alpha = 0.1$) at 10,000 and 16,000 steps. For the WiCLP, a window size of 5 was used. All network parameters were initialized to zero, leveraging the skip connection to ensure that the initial output 689 matched the CLIP text embeddings. Our implementation is based on the Diffusers² library, utilizing their 690 modules, models, and checkpoints to build and train our models. This comprehensive setup ensured that our 691 method was rigorously tested under controlled conditions, providing a robust evaluation of its performance. 692

D.2 EXTENDED EVALUATION

Human Evaluation We conducted a human evaluation in which participants compared images generated by SD v1.4 and SD v1.4 + WiCLP, selecting the image that best matched the given prompt (Figure 18). Five evaluators were presented with 200 randomly selected image pairs, evaluating a total of 1000 image-caption pairs.

TIFA Metric. To provide a more comprehensive evaluation, in addition to the disentangled BLIP-VQA score proposed by Huang et al. (2023), we also incorporate the TIFA metric (Hu et al., 2023). TIFA (Text-to-Image Faithfulness Evaluation with Question Answering) is an automated evaluation method that measures how faithfully a generated image corresponds to its textual input via visual question answering (VQA). It generates

^{704 &}lt;sup>2</sup>https://github.com/huggingface/diffusers



Zero-shot Attention Reweighting Original Text-embeddings

Figure 11: Visualization of some images generated with same set of seeds using original text-embeddings of prompt "a blue car and a brown cow" and text-embeddings that are obtained as the result of zero-shot reweighting of attention matrix.

multiple question-answer pairs from the text input using a language model, then evaluates image faithfulness by determining whether existing VQA models can accurately answer these questions based on the image. As a reference-free metric, TIFA offers fine-grained and interpretable assessments of image quality.

Using TIFA, we observed that SD v1.4 and SD v2 achieved scores of 0.6598 and 0.7735, respectively. Notably, the scores for WiCLP applied on top of SD v1.4 and SD v2 improved to 0.7462 and 0.8133, respectively, demonstrating the enhanced performance of our approach.

FID Score Comparison Our method results in a modest increase in FID score on MS-COCO prompts compared to the base models, as shown in Table 3. However, this increase is less pronounced than in other baselines—for example, SD v2 + WiCLP scores 27.40, whereas GORS reaches 30.54.

D.3 CLP AND WICLP VISUALIZATION

In this section, we provide additional visualizations comparing CLP, WiCLP, and baseline models in Figures 14, 15.

- VISUALIZATION OF CROSS-ATTENTIONS 749 D.4
- 751 In this section, we provide additional cross-attention map visualizations in Figures 14 and 15.

752 D.5 VISUALIZATION OF SWITCH-OFF

In this section, we present more qualitative samples illustrating the effect of SWITCH-OFF at different timestep thresholds for various prompts in Figures 16 and 17.

D.6 CHOICE OF WINDOW LENGTH IN WICLP

One might suggest that instead of using token-wise linear projection (CLP) or a window-based linear projection with a limited window (WiCLP), employing a linear projection that considers all tokens when finding a better embedding for each token might yield better results. However, our thorough quantitative study and experiments tested various window sizes for WiCLP. We found that using a window size of 5 achieves the highest performance.

799		SD v1.4	CLP	SD v2	WiCLP
800					
801					1000
802	A blue bowl and				
803	a red train				
804					7
805				All a Balan	
806					
807	A 1-1			1 1 1 1 1	
808	A blue bench and	Carlo Carlo			
809	a green bowr				
810			A Barran		
811		-			
812			~		ONESNER A PROCEED
813	A blue backpack and				
814	a red book				
815					
816					
817					the second se
818				100	-
819	A black and white cat		100		-
820	sitting in a green bowl				
821					
822				For show we want	
823			4		
824	A brown boat and	and the second second			
825	a blue cat			20 500 12	
826	a blue eat				
021					
020 920				A () A	
920				50. XI	
831	A brown book and		-		
832	a red sheep	· Como			
833			Winds.		
834		II.			
835					
836			SIC		
837	A fluffy towel and	~			
838	a glass cup				
839			T	-	
840					
841					
842	A plastic container and			19	
843	a fluffy teddy bear				ALA
844					
845			Company of the local sector		

Figure 12: Caption

846		SD v1.4	CLP	SD v2	WiCLP
847					
848					
849	A red apple and				•
850	a green train				
851			Real Property of the second se		
852				And the second	
853					
854		And			
855	A red chair and				
856	a gold clock				
857		7 7			
858					
859					763
860	A red pen and				
861	a blue notebook				
862					
863					and the state of t
864					
865					
866	A round cookie and				
867	a square container	(CA)	1. C. C. P.		
868		NED D	WE W		
869					
870					A CONTRACT OF ANY ANY
871			and and the		
872	A wooden floor and	AT -	1		
873	a nuny rug				and the second second
874				anter.	
875					
876					
877	The leather jacket and fluffy				Constants !!
8/8	scarf keep the cold at bay				
8/9					
000					
001					
002		(lille			
003	Wooden pencil and a glass plate				
004					
200					
000					
00/		19. J. C. S.			
000			CACA!		
800	A green leaf and a vellow butterfly	· · · ·			warman .
801	a yenow butterny	¥ (\sim	
800		1. 1. 1.			
002		And the second se	the second s		and the second

Figure 13: Caption







Figure 16: Qualitative results showing the Ampact of SWITCH-OFF with varying thresholds T



Figure 17: Qualitative results showing the Ampact of SWITCH-OFF with varying thresholds T

