

# DIRECTIONAL TEXTUAL INVERSION FOR PERSONALIZED TEXT-TO-IMAGE GENERATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Textual Inversion (TI) is efficient for text-to-image personalization but often fails on complex prompts. We identify embedding norm inflation as a key cause and show that token semantics are primarily encoded by embedding direction. We propose *Directional Textual Inversion* (DTI), which fixes the embedding magnitude to an in-distribution scale and optimizes only direction with a simple MAP objective using a von Mises-Fisher prior. DTI improves prompt fidelity over existing embedding optimization baselines while maintaining competitive subject similarity. Furthermore, we demonstrate its ease of integration and creative applications.

## 1 INTRODUCTION

Personalization in text-to-image generation adapts models to user-specific concepts. While parameter fine-tuning methods like DreamBooth (Ruiz et al., 2023) are effective, they are computationally intensive. In contrast, embedding optimization methods, such as Textual Inversion (TI) (Gal et al., 2023a), offer a lightweight alternative by updating only token embeddings. This efficiency has made TI a foundational component in numerous personalization frameworks (Hao et al., 2023; Kumari et al., 2023; Tewel et al., 2023b; Lee et al., 2024).

Despite its utility, TI faces critical limitations. Constraining complex visual concepts to a single vector often compromises prompt fidelity and necessitates lengthy optimization. Recent attempts to enrich the embedding space (Voynov et al., 2023; Alaluf et al., 2023) introduce computational overhead without addressing the underlying optimization dynamics that govern semantic alignment.

This paper presents a systematic analysis of TI, revealing that semantic information is predominantly encoded in the embedding direction. We demonstrate, both theoretically and empirically, that unbounded embedding magnitudes during optimization are a primary source of instability, significantly impairing image-text alignment.

Building on these insights, we introduce **Directional Textual Inversion (DTI)**. Unlike standard TI, DTI decouples magnitude from direction; it maintains the magnitude consistent with pre-trained distributions while focusing optimization exclusively on direction. We formulate this as a Maximum a Posteriori (MAP) estimation problem, utilizing a von Mises-Fisher (vMF) prior to regularize the hyperspherical latent space. This approach preserves the efficiency of TI while ensuring robust semantic alignment.

Our evaluation demonstrates that DTI consistently outperforms conventional TI and enhancements like CrossInit (Pang et al., 2024a), achieving superior semantic fidelity with greater efficiency. Furthermore, our directionally optimized embeddings enable novel applications, such as smooth interpolation between personalized concepts.

## 2 ANALYZING TOKEN EMBEDDING GEOMETRY

### 2.1 EMPIRICAL MOTIVATION: DIRECTION ENCODES SEMANTICS

Aligning with previous studies (Mikolov et al., 2013; Pennington et al., 2014), we find that token embeddings encode meaning primarily in direction, not magnitude. We validate this by comparing nearest neighbors: cosine similarity (direction-sensitive) yields semantically coherent neighbors for ‘apple’ (*top-5: apples, fruit, peach, pear, egg*), whereas Euclidean distance (magnitude-sensitive)

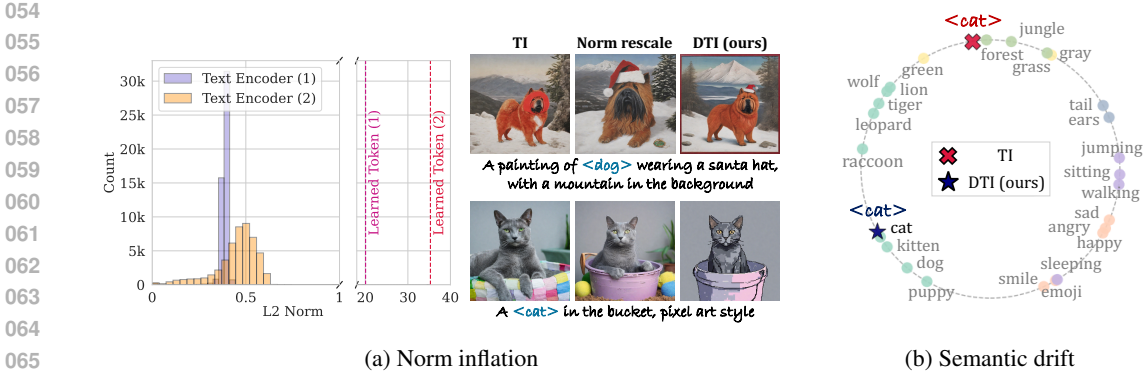


Figure 1: **Empirical Motivation.** Standard TI degrades prompt fidelity due to (a) **excessive norms** compared to the original vocabulary and (b) **semantic drift** away from related concepts. DTI addresses these by preserving both norm scale and directional integrity.

finds unrelated tokens (*top-5*:  $U+2069$ , *altrin*, *lestwe*, *heartnews*, *samanthaprabhu*). Figure 1b shows that standard TI often neglects this, causing semantic drift where learned embeddings (e.g.,  $\langle \text{cat} \rangle$ ) move away from related concepts. This motivates our directional optimization approach.

## 2.2 WHY LARGE MAGNITUDES LEAD TO LOW TEXT FIDELITY

**Effect I: Positional information is attenuated (see Lemma 1).** After LayerNorm/RMSNorm layer, the normalized signal that feeds attention/MLP becomes less sensitive to small additive terms as  $m$  grows. Positional information contributes  $\mathcal{O}(1/m)$  to the normalized signal  $\text{Norm}(mv + p)$ . Intuitively, a very large-norm token *forgets where it is in the sequence*, weakening contextualization, resulting in omission of details such as style and background (see Figure 1).

**Effect II: Residual updates stagnate (see Lemma 2).** The residual updates,  $F_\ell(\text{Norm}(\mathbf{x}^{(\ell)}))$ , are computed from a *normalized* inputs and thus have a bounded magnitude. When this bounded update is added through the skip connection to a large vector  $\mathbf{x}^{(\ell)}$ , the *relative change* (i.e., turning angle of the hidden state’s direction) becomes tiny, decreasing in proportion to  $1/\|\mathbf{x}^{(\ell)}\|$ . In other words, large-norm hidden states become *stuck* in their direction and are difficult for subsequent layers to refine. This *residual stagnation* accumulates across layers, severely limiting the total directional change the initial token can undergo, as formalized in the following proposition and corollary.

**Proposition 1** (Accumulated directional drift across  $L$  pre-norm blocks). *Let  $\mathbf{x}^{(0)} \neq \mathbf{0}$  and  $\mathbf{x}^{(\ell+1)} = \mathbf{x}^{(\ell)} + F_\ell(\text{Norm}(\mathbf{x}^{(\ell)}))$  for  $\ell = 0, \dots, L - 1$ . Let  $B_\ell := \sup_{\mathbf{u} \in S} \|F_\ell(\mathbf{u})\|_2 < \infty$ , and  $S_L := \sum_{j=0}^{L-1} B_j$ . Assume  $\|\mathbf{x}^{(0)}\|_2 > S_L$ , then*

$$\angle(\mathbf{x}^{(0)}, \mathbf{x}^{(L)}) \leq \frac{\pi}{2} \sum_{\ell=0}^{L-1} \frac{B_\ell}{\|\mathbf{x}^{(0)}\|_2 - \sum_{j<\ell} B_j} \leq \frac{\pi}{2} \frac{S_L}{\|\mathbf{x}^{(0)}\|_2 - S_L}.$$

**Corollary 1** (Scaling  $\Rightarrow$  directional freezing). *With the notation of Proposition 1, for any  $\alpha > 1$ ,*

$$\angle(\alpha \mathbf{x}^{(0)}, \mathbf{x}^{(L)}(\alpha)) \leq \frac{\pi}{2} \frac{S_L}{\alpha \|\mathbf{x}^{(0)}\| - S_L} \xrightarrow{\alpha \rightarrow \infty} 0,$$

where  $\mathbf{x}^{(L)}(\alpha)$  denotes the depth- $L$  output when the initial token is  $\alpha \mathbf{x}^{(0)}$ .

Together, this explains TI’s reduced text fidelity: excessive token magnitude inhibits contextual integration, causing the personalized subject to overshadow prompt constraints such as style and background. This necessitates the explicit magnitude control framework introduced in the subsequent section. We provide formal proofs in Appendix C and empirical validations in Appendix D.

## 3 METHOD: DIRECTIONAL TEXTUAL INVERSION

Based on our observation and analysis on previous section that token embeddings exhibit strong directional characteristics, we introduce *Directional Textual Inversion* (DTI), a framework that

108 optimizes an embedding’s direction with in-distribution norm to enhance text fidelity in personalized  
 109 text-to-image generation.

### 111 3.1 OPTIMIZING ONLY DIRECTION ON THE HYPERSPHERE

112 We reformulate TI by decoupling the magnitude and direction of the learnable token embedding  
 113  $e \in \mathbb{R}^d$ . The embedding can be expressed as:

$$114 \quad e = m^* \mathbf{v}, \quad \mathbf{v} \in \mathbb{S}^{d-1}. \quad (1)$$

115 Here,  $\mathbb{S}^{d-1} = \{\mathbf{u} \in \mathbb{R}^d : \|\mathbf{u}\|_2 = 1\}$  denotes the unit sphere. We fix the magnitude  $m^*$  and optimize  
 116 only the direction ( $\mathbf{v}$ ). Specifically, we set  $m^*$  to be an *in-distribution* magnitude derived from the  
 117 frozen vocabulary of text encoder (e.g., the average norm). In this way, optimization focuses on  
 118 semantic information in direction while avoiding out-of-distribution (OOD) norms.

119 Since the parameter space is the unit sphere, Euclidean updates drift off-manifold, making  
 120 AdamW (Loshchilov & Hutter, 2017) (default optimizer used in TI-like methods) not suitable.  
 121 To solve this, we use Riemannian stochastic gradient descent (RSGD) (Bonnabel, 2013) with tangent-  
 122 space projection and retraction. See Algorithm 1 and Appendix E.1 for further details.

### 125 3.2 MAXIMUM A POSTERIORI FORMULATION WITH A DIRECTIONAL vMF PRIOR

126 To incorporate directional prior, we formulate the optimization for the optimal direction  $\mathbf{v}^*$  as a  
 127 Maximum A Posteriori (MAP) estimation problem. Given a dataset of images  $\mathcal{D} = \{z_1, \dots, z_n\}$ ,  
 128 the MAP estimate is found by maximizing the posterior probability:

$$129 \quad \mathbf{v}^* = \arg \max_{\mathbf{v}} p(\mathbf{v} | \mathcal{D}) \propto \arg \max_{\mathbf{v}} [\log p(\mathcal{D} | \mathbf{v}) + \log p(\mathbf{v})]. \quad (2)$$

130 Minimizing the negative log-posterior is equivalent to minimizing a loss function composed of a data  
 131 term and a prior term:  $\mathcal{L}(\mathbf{v}) = \mathcal{L}_{\text{data}}(\mathbf{v}) + \mathcal{L}_{\text{prior}}(\mathbf{v})$ .

132 The data term,  $\mathcal{L}_{\text{data}} = -\log p(\mathcal{D} | \mathbf{v})$ , is the negative log-likelihood of the images given the  
 133 direction. Following standard practice for diffusion models (Ho et al., 2020), we use the mean  
 134 squared error (MSE) between the true and predicted noise as the objective:  $\mathcal{L}_{\text{data}}(\mathbf{v}) := \mathbb{E}_{z,t,\epsilon,c} [\|\epsilon -$   
 135  $\epsilon_\theta(z_t, t, c(\mathbf{v}))\|_2^2]$ . Here,  $\epsilon_\theta$  and  $c(\cdot)$  are the diffusion model and text encoder, respectively. The  
 136 Euclidean gradient of this objective,  $\mathbf{g}_{\text{euc}} = \nabla_{\mathbf{v}} \mathcal{L}$ , is used in the RSGD update.

137 For the prior term,  $-\log p(\mathbf{v})$ , we use a von Mises-Fisher (vMF) distribution on the direction  $\mathbf{v}$   
 138 (detailed justification in Appendix E.2). The vMF distribution is a probability distribution on the  
 139  $(d-1)$ -sphere, analogous to the Gaussian distribution in Euclidean space. It is parameterized by a  
 140 mean direction  $\boldsymbol{\mu} \in \mathbb{S}^{d-1}$  and a concentration parameter  $\kappa \geq 0$ . The probability density function is  
 141 given by:

$$142 \quad p(\mathbf{v} | \boldsymbol{\mu}, \kappa) = \frac{\kappa^{d/2-1}}{(2\pi)^{d/2} I_{d/2-1}(\kappa)} \exp(\kappa \boldsymbol{\mu}^\top \mathbf{v}), \quad (3)$$

143 where  $I_{d/2-1}$  is the modified Bessel function of the first kind. Here, we work with unnormalized  
 144 density:  $p(\mathbf{v}) \propto \exp(\kappa \boldsymbol{\mu}^\top \mathbf{v})$ . Ignoring constants, the negative log-prior yields our regularization  
 145 term,  $\mathcal{L}_{\text{prior}}(\mathbf{v}) = -\kappa \boldsymbol{\mu}^\top \mathbf{v}$ .

146 **Selection of vMF parameters.** We set mean direction  $\boldsymbol{\mu}$  to the constant, normalized embedding of  
 147 the corresponding class token (e.g., ‘dog’). Since estimating concentration parameter  $\kappa$  is non-trivial,  
 148 we treat it as a hyperparameter that controls the strength of the prior. Based on a grid search (optimal  
 149 range  $5e-5$  to  $2e-4$ ), we fixed  $\kappa$  to  $1e-4$  for all experiments. See Appendix F.2 and F.4 for details.

## 153 4 EXPERIMENTS

### 154 4.1 EXPERIMENTAL SETUPS

155 We evaluated our method using all reference images from the DreamBooth dataset (Ruiz et al.,  
 156 2023) across 40 prompts, including 10 complex additions. Our primary model was Stable Diffusion  
 157 XL (SDXL) (Podell et al., 2024), with additional experiments on SANA 1.5 (Xie et al., 2025)  
 158 to demonstrate applicability to recent architectures. We compared DTI primarily against Textual  
 159 Inversion (TI) (Gal et al., 2023a) and CrossInit (Pang et al., 2024a), deferring other baselines to  
 160 Appendix F.3. Following established protocols (Ruiz et al., 2023; Kumari et al., 2023), we measured  
 161 subject fidelity using DINOv2 (Oquab et al., 2023) cosine similarity and image-text alignment using  
 SigLIP (Zhai et al., 2023). Implementation details are in Appendix F.1.

Table 1: Our DTI consistently improves baselines by generating outputs with enhanced text fidelity while maintaining subject similarity.

Methods	SDXL		SANA 1.5-1.6B		SANA 1.5-4.8B	
	Image	Text	Image	Text	Image	Text
TI	<b>0.561</b>	0.292	<b>0.480</b>	<u>0.621</u>	<u>0.446</u>	<u>0.646</u>
TI-rescaled	0.243	0.466	0.253	0.655	0.287	0.548
CrossInit	0.545	0.464	0.344	0.614	0.299	0.622
<b>DTI (ours)</b>	0.450	<b>0.522</b>	<u>0.479</u>	<b>0.744</b>	<b>0.452</b>	<b>0.757</b>

Table 2: Ablation studies. We tested and confirmed the effectiveness of every component of our DTI.

Optimizer	$m^*$	$\kappa \times 10^{-3}$	Image	Text
AdamW	mean	0.1	0.335	0.463
RSGD	min	0.1	0.030	0.074
RSGD	5.0 (OOD)	0.1	0.383	0.373
RSGD	mean	0.0	<b>0.507</b>	<u>0.436</u>
RSGD	mean	0.5	0.278	<b>0.688</b>
RSGD	mean	0.1	<u>0.450</u>	<u>0.522</u>



Figure 2: **Subject interpolation.** DTI enables smooth, coherent conceptual transitions, expanding creative personalization capabilities.

## 4.2 MAIN RESULTS

We present the main results in the main paper, while more experiments including LoRA fine-tuning, stylization, my object in my style, and face personalization can be found in Appendix F.

**Quantitative results.** In Table 1, we quantitatively evaluate DTI along two axes: subject similarity and text-prompt fidelity. Overall, the results clearly demonstrate the advantage of DTI over the baselines. Additional comparisons with further baselines on other Stable Diffusion variants are provided in Appendix F.3.

**Ablation study.** We performed ablation study to verify the effectiveness of components of our DTI, including the optimization space, the embedding magnitude  $m$ , and the concentration parameter of vMF distribution  $\kappa$ . The results are summarized in Table 2, with further analyses in Appendix F.4.

**Human evaluation.** To further examine the effectiveness of our method, we conducted a large scale user study (100 participants via *Amazon Mechanical Turk*) to measure real-world user preferences, with the details in Appendix F.5.

**Embedding interpolation for creative applications.** We demonstrate DTI’s creative potential through embedding interpolation experiments (Figure 2). Leveraging its unit-spherical embedding space, DTI enables coherent transitions via spherical linear interpolation (SLERP), unlike the linear interpolation used by TI which often fails. DTI seamlessly blends distinct concepts, such as merging a dog and a teapot into a hybrid object, or smoothly transitioning between animals like a dog and a cat. It also excels in nuanced face personalization, plausibly interpolating between a young boy and an older woman while maintaining coherence. These results highlight DTI as a powerful tool for intuitive concept blending and creative applications. Additional results are in Appendix F.7.

## 5 DISCUSSION & CONCLUSION

As DTI focuses on text prompt fidelity, for applications requiring higher subject fidelity, we recommend integrating DTI with lightweight fine-tuning methods. Empirically, we show our DTI effectively complements lightweight fine-tuning methods like LoRA in Figure 5 and 10.

In this paper, we identified embedding norm inflation as a critical bottleneck for text-to-image alignment. By leveraging the underexplored directional characteristics of the token embedding space, we introduced DTI which optimizes embedding direction while constraining the norm to an in-distribution scale via a novel directional prior. DTI demonstrably enhances prompt fidelity, offering a robust solution for precise and controllable personalization in generative AI.

## REFERENCES

- 216  
217  
218 Yuval Alaluf, Elad Richardson, Gal Metzer, and Daniel Cohen-Or. A neural space-time representation  
219 for text-to-image personalization. *ACM Transactions on Graphics (TOG)*, 42(6):1–10, 2023.
- 220 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint*  
221 *arXiv:1607.06450*, 2016.
- 222  
223 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang,  
224 Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*,  
225 2025.
- 226 Silvere Bonnabel. Stochastic gradient descent on Riemannian manifolds. *IEEE Transactions on*  
227 *Automatic Control*, 58(9):2217–2229, September 2013. ISSN 0018-9286, 1558-2523.
- 228  
229 Hong Chen, Yipeng Zhang, Simin Wu, Xin Wang, Xuguang Duan, Yuwei Zhou, and Wenwu Zhu.  
230 Disenbooth: Identity-preserving disentangled tuning for subject-driven text-to-image generation.  
231 In *International Conference on Learning Representations*, 2023a.
- 232  
233 Wenhui Chen, Hexiang Hu, Yandong Li, Nataniel Ruiz, Xuhui Jia, Ming-Wei Chang, and William W  
234 Cohen. Subject-driven text-to-image generation via apprenticeship learning. In *Advances in Neural*  
235 *Information Processing Systems*, 2023b.
- 236 Minhyung Cho and Jaehyung Lee. Riemannian approach to batch normalization. In *Advances in*  
237 *Neural Information Processing Systems*, October 2017.
- 238  
239 John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and  
240 stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- 241 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel  
242 Cohen-Or. An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual  
243 Inversion. In *International Conference on Learning Representations*, 2023a.
- 244 Rinon Gal, Moab Arar, Yuval Atzmon, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or.  
245 Encoder-based domain tuning for fast personalization of text-to-image models. *ACM Transactions*  
246 *on Graphics (TOG)*, 42(4):1–13, 2023b.
- 247  
248 Yuchao Gu, Xintao Wang, Jay Zhangjie Wu, Yujun Shi, Yunpeng Chen, Zihan Fan, Wuyou Xiao,  
249 Rui Zhao, Shuning Chang, Weijia Wu, et al. Mix-of-show: Decentralized low-rank adaptation for  
250 multi-concept customization of diffusion models. In *Advances in Neural Information Processing*  
251 *Systems*, 2023.
- 252 Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang. Svdiff:  
253 Compact parameter space for diffusion fine-tuning. In *Proceedings of the IEEE/CVF Conference*  
254 *on Computer Vision and Pattern Recognition*, pp. 7323–7334, 2023.
- 255  
256 Shaozhe Hao, Kai Han, Shihao Zhao, and Kwan-Yee K Wong. ViCo: Plug-and-play visual condition  
257 for personalized text-to-image generation, 2023.
- 258 Geoffrey E Hinton, Nitish Srivastava, and Kevin Swersky. Lecture 6e: Rmsprop. Coursera, Nonlinear  
259 aural component analysis of the cochlea, 2012. URL [http://www.cs.toronto.edu/  
260 ~tijmen/csc321/slides/lecture\\_slides\\_lec6.pdf](http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf).
- 261  
262 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. In *Advances*  
263 *in Neural Information Processing Systems*, 2020.
- 264 Shoaib Jameel and Steven Schockaert. Word and document embedding with vmf-mixture priors on  
265 context word vectors. In *Proceedings of the Annual Meeting of the Association for Computational*  
266 *Linguistics*, 2019.
- 267  
268 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative  
269 adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
*Pattern Recognition*, pp. 4401–4410, 2019.

- 270 Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *International*  
271 *Conference on Learning Representations*, 2015.
- 272
- 273 Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-Concept  
274 Customization of Text-to-Image Diffusion. In *Proceedings of the IEEE/CVF Conference on*  
275 *Computer Vision and Pattern Recognition*, June 2023.
- 276 Kyungmin Lee, Sangkyung Kwak, Kihyuk Sohn, and Jinwoo Shin. Direct Consistency Optimization  
277 for Compositional Text-to-Image Personalization. In *Advances in Neural Information Processing*  
278 *Systems*, 2024.
- 279
- 280 Dongxu Li, Junnan Li, and Steven Hoi. Blip-diffusion: Pre-trained subject representation for  
281 controllable text-to-image generation and editing. In *Advances in Neural Information Processing*  
282 *Systems*, 2023.
- 283 Xianming Li and Jing Li. Aoe: Angle-optimized embeddings for semantic textual similarity. In  
284 *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume*  
285 *1: Long Papers)*, pp. 1825–1839, 2024.
- 286 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*  
287 *arXiv:1711.05101*, 2017.
- 288
- 289 Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *International*  
290 *Conference on Learning Representations*, 2019.
- 291
- 292 Jian Ma, Junhao Liang, Chen Chen, and Haonan Lu. Subject-diffusion: Open domain personalized  
293 text-to-image generation without test-time fine-tuning. In *ACM SIGGRAPH 2024 Conference*  
294 *Papers*, pp. 1–12, 2024.
- 295 Yu Meng, Jiaxin Huang, Guangyuan Wang, Chao Zhang, Honglei Zhuang, Lance Kaplan, and Jiawei  
296 Han. Spherical text embedding. In *Advances in Neural Information Processing Systems*, volume 32,  
297 2019.
- 298
- 299 Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations  
300 of words and phrases and their compositionality. In *Advances in Neural Information Processing*  
301 *Systems*, volume 26, 2013.
- 302 Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob  
303 Mcgrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards Photorealistic Image Generation and  
304 Editing with Text-Guided Diffusion Models. In *Proceedings of the International Conference on*  
305 *Machine Learning*, pp. 16784–16804. PMLR, 2022.
- 306
- 307 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov,  
308 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas  
309 Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael  
310 Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut,  
311 Armand Joulin, and Piotr Bojanowski. DINOv2: Learning Robust Visual Features without  
Supervision. *Transactions on Machine Learning Research*, July 2023. ISSN 2835-8856.
- 312
- 313 Lianyu Pang, Jian Yin, Haoran Xie, Qiping Wang, Qing Li, and Xudong Mao. Cross initialization  
314 for face personalization of text-to-image models. In *Proceedings of the IEEE/CVF Conference on*  
315 *Computer Vision and Pattern Recognition*, 2024a.
- 316 Lianyu Pang, Jian Yin, Baoquan Zhao, Feize Wu, Fu Lee Wang, Qing Li, and Xudong Mao.  
317 Attdreambooth: Towards text-aligned personalized text-to-image generation. In *Advances in*  
318 *Neural Information Processing Systems*, 2024b.
- 319
- 320 NaHyeon Park, Kunhee Kim, and Hyunjung Shim. TextBoost: Towards One-Shot Personalization of  
321 Text-to-Image Models via Fine-tuning Text Encoder. *arXiv preprint arXiv:2409.08248*, 2024.
- 322
- 323 Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word  
representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language*  
*Processing*, pp. 1532–1543, 2014.

- 324 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
325 Penna, and Robin Rombach. SDXL: Improving Latent Diffusion Models for High-Resolution  
326 Image Synthesis. In *International Conference on Learning Representations*, 2024.  
327
- 328 Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian  
329 Weller, and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetuning. In  
330 *Advances in Neural Information Processing Systems*, 2023.
- 331 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen,  
332 and Ilya Sutskever. Zero-Shot Text-to-Image Generation. In *Proceedings of the International  
333 Conference on Machine Learning*, 2021.  
334
- 335 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical Text-  
336 Conditional Image Generation with CLIP Latents, April 2022.
- 337 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
338 Resolution Image Synthesis with Latent Diffusion Models. In *Proceedings of the IEEE/CVF  
339 Conference on Computer Vision and Pattern Recognition*, 2022.  
340
- 341 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.  
342 DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation. In  
343 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- 344 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Wei Wei, Tingbo Hou, Yael Pritch, Neal Wadhwa,  
345 Michael Rubinstein, and Kfir Aberman. Hyperdreambooth: Hypernetworks for fast personalization  
346 of text-to-image models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
347 Pattern Recognition*, pp. 6527–6536, 2024.  
348
- 349 Yoad Tewel, Rinon Gal, Gal Chechik, and Yuval Atzmon. Key-locked rank one editing for text-to-  
350 image personalization. In *ACM SIGGRAPH 2023 conference proceedings*, pp. 1–11, 2023a.
- 351 Yoad Tewel, Rinon Gal, Gal Chechik, and Yuval Atzmon. Key-Locked Rank One Editing for Text-  
352 to-Image Personalization. In *Proc. SIGGRAPH, SIGGRAPH '23*, New York, NY, USA, 2023b.  
353 Association for Computing Machinery. ISBN 979-8-4007-0159-7.  
354
- 355 Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman. P+: Extended Textual Condi-  
356 tioning in Text-to-Image Generation, March 2023.
- 357 Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through align-  
358 ment and uniformity on the hypersphere. In *Proceedings of the International Conference on  
359 Machine Learning*, pp. 9929–9939. PMLR, 2020.  
360
- 361 Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. Elite: Encoding  
362 visual concepts into textual embeddings for customized text-to-image generation. In *Proceedings  
363 of the IEEE/CVF International Conference on Computer Vision*, pp. 15943–15953, 2023.
- 364 Feize Wu, Yun Pang, Junyi Zhang, Lianyu Pang, Jian Yin, Baoquan Zhao, Qing Li, and Xudong  
365 Mao. Core: Context-regularized text embedding learning for text-to-image personalization. In  
366 *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025.  
367
- 368 Enze Xie, Junsong Chen, Junyu Chen, Han Cai, Haotian Tang, Yujun Lin, Zhekai Zhang, Muyang Li,  
369 Ligeng Zhu, Yao Lu, et al. Sana: Efficient high-resolution image synthesis with linear diffusion  
370 transformers. In *International Conference on Learning Representations*, 2025.
- 371 Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt  
372 adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023.  
373
- 374 Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan,  
375 Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin  
376 Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling Autoregressive Models for Content-  
377 Rich Text-to-Image Generation. *Transactions on Machine Learning Research*, August 2022. ISSN  
2835-8856.

378 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language  
379 image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,  
380 pp. 11975–11986, 2023.

381 Biao Zhang and Rico Sennrich. Root Mean Square Layer Normalization. In *Advances in Neural  
382 Information Processing Systems*, 2019.

383  
384 Dingyi Zhang, Yingming Li, and Zhongfei Zhang. Deep metric learning with spherical embedding.  
385 In *Advances in Neural Information Processing Systems*, 2020.

386  
387 Xulu Zhang, Xiao-Yong Wei, Jinlin Wu, Tianyi Zhang, Zhaoxiang Zhang, Zhen Lei, and Qing Li.  
388 Compositional inversion for stable diffusion models. In *Proceedings of the AAAI Conference on  
389 Artificial Intelligence*, 2024a.

390 Yanbing Zhang, Mengping Yang, Qin Zhou, and Zhe Wang. Attention calibration for disentangled  
391 text-to-image personalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
392 and Pattern Recognition*, pp. 4764–4774, 2024b.

393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

## 432 A RELATED WORK

### 433 A.1 PERSONALIZED TEXT-TO-IMAGE GENERATION

434 Recent advancements in text-to-image (T2I) generation have considerably expanded the creative  
 435 capabilities and flexibility of generative models (Ramesh et al., 2021; Rombach et al., 2022; Nichol  
 436 et al., 2022; Ramesh et al., 2022; Yu et al., 2022; Podell et al., 2024). Despite these innovations,  
 437 natural language inherently struggles to precisely convey nuanced, user-specific concepts. This  
 438 inherent limitation has driven the development of personalization methods, which allow users to  
 439 generate images reflecting unique concepts with creative prompts.  
 440

441 Textual Inversion (Gal et al., 2023a), which is most well-known for its lightweight integration to  
 442 many other personalization works, uses embedding optimization by introducing learnable tokens for  
 443 personalized information without model modification. Subsequent work explored diverse embedding  
 444 strategies (Voynov et al., 2023; Alaluf et al., 2023; Wu et al., 2025; Zhang et al., 2024a), often with  
 445 demanding excessive computational costs. Among them, CrossInit (Pang et al., 2024a) offered an  
 446 efficient initialization strategy with minimal overhead, replacing initialization tokens with the output  
 447 of text encoder and using regularization loss.  
 448

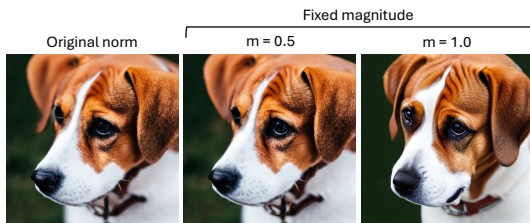
449 In contrast, fine-tuning based methods such as DreamBooth (Ruiz et al., 2023) achieve high sub-  
 450 ject fidelity, but require significant computational resources compared to embedding optimization  
 451 methods (Kumari et al., 2023; Han et al., 2023; Gu et al., 2023; Chen et al., 2023a; Tewel et al.,  
 452 2023a; Zhang et al., 2024b; Qiu et al., 2023; Pang et al., 2024b). More recently, Park et al. (2024)  
 453 proposed fine-tuning text encoder instead of image generator for efficiency, but they still demand  
 454 more parameters compared to embedding optimization methods.

455 Meanwhile, there exists a line of encoder-based approaches (Wei et al., 2023; Ruiz et al., 2024; Ye  
 456 et al., 2023; Gal et al., 2023b; Chen et al., 2023b; Li et al., 2023; Pang et al., 2024b; Ma et al., 2024)  
 457 that offer fast inference, but they necessitate substantial pre-training.  
 458

### 459 A.2 DIRECTIONAL EMBEDDING SPACE

460 A number of prior works has emphasized constraining embedding representations to the hypersphere.  
 461 These include using vMF mixtures for directional clustering (Jameel & Schockaert, 2019), normalizing  
 462 norms for face recognition (Meng et al., 2019), angle-optimized embeddings to address cosine  
 463 saturation (Li & Li, 2024), and spherical constraints for uniform document clustering (Zhang et al.,  
 464 2020). Wang & Isola (2020) offered theoretical support for hyperspherical constraints in contrastive  
 465 learning. Our method aligns with this trend by modeling embeddings as directional distributions but  
 466 uniquely decomposes and explicitly optimizes textual embedding direction using a vMF prior within  
 467 Textual Inversion framework.  
 468

## 469 B EMBEDDING NORM AND DIRECTION



481 Figure 3: **Effect of magnitude change.** We set the magnitude to a fixed value to analyze the impact  
 482 of magnitude changes. The resulting outputs show no noticeable difference.

483 We altered the magnitude of the token as exemplified in Figure 3. However, the resulting output re-  
 484 mained mostly unchanged. This indicates that minor adjustments to the magnitude do not significantly  
 485 affect the outcome.

Table 3: **Nearest tokens under different measures.** We show the nearest tokens to the query words ‘study’ and ‘writing’ using both cosine similarity and Euclidean distance.

Query	Cosine	Euclidean
study	studies, studying, research, bookclub, reading, studied, sketches, measurements, thumbnail	U+3160, texanscheer, asober, instaweatherpro, mydayin, premiosmtvmiaw, tairp, thepersonalnetwork, U+2412
writing	writer, write, written, writ, writers, writings, recording, blogging, wrote	phdlife, poetryday, tomorrowspaper, urstrulymahesh, @_--, twitterkurds, asober, fakespeare, jamiedor

In Table 3, we provide additional examples illustrating the nearest words retrieved for each query under different similarity measures, which strongly correlate with either direction or magnitude. Our analysis reveals that cosine similarity retrieves words that share semantic meaning with the query. Conversely, Euclidean distance is significantly affected by embedding magnitude, often retrieving words with limited or no semantic relevance. This demonstrates that semantic meaning is predominantly associated with embedding direction rather than magnitude. Note that words beginning with U+ denote Unicode.

## C PROOFS FOR THEORETICAL STATEMENTS

### C.1 SETUP

**Pre-norm block.** We study *pre-norm* Transformer blocks

$$\mathbf{x}^{(\ell+1)} = \mathbf{x}^{(\ell)} + F_\ell(\text{Norm}(\mathbf{x}^{(\ell)})), \quad \ell = 0, \dots, L - 1, \quad (4)$$

where  $\text{Norm} \in \{\text{LayerNorm}, \text{RMSNorm}\}$  (with optional affine  $(\gamma, \beta)$  absorbed into  $F_\ell$ ).

**Scale invariance.** For normalizations, we use the standard, scale-invariant definitions:

$$\text{RMSN}(\mathbf{x}) = \sqrt{d} \frac{\mathbf{x}}{\|\mathbf{x}\|_2}, \quad \text{LN}(\mathbf{x}) = \sqrt{d} \frac{\mathbf{C}\mathbf{x}}{\|\mathbf{C}\mathbf{x}\|_2}, \quad \mathbf{C} := \mathbf{I} - \frac{1}{d} \mathbf{1}\mathbf{1}^\top. \quad (5)$$

Thus  $\text{RMSN}(s\mathbf{x}) = \text{RMSN}(\mathbf{x})$  and  $\text{LN}(s\mathbf{x}) = \text{LN}(\mathbf{x})$  for all  $s > 0$ . Please refer to original papers (Ba et al., 2016; Zhang & Sennrich, 2019) for further details.

**Token decomposition.** For the input token, we denote  $\mathbf{x}^{(0)} = m\mathbf{v} + \mathbf{p}$  with  $m > 0$ ,  $\|\mathbf{v}\|_2 = 1$ , and (optional) absolute positional embedding  $\mathbf{p} \in \mathbb{R}^d$ .

**Bounded sub-layers.** Define  $\mathcal{S} = \{\text{Norm}(\mathbf{z}) : \mathbf{z} \neq \mathbf{0}\}$ . Since  $\text{Norm}$  maps into a fixed scale, bounded set and  $F_\ell$  (attention/MLP plus projections) is continuous on bounded sets,

$$B_\ell := \sup_{\mathbf{u} \in \mathcal{S}} \|F_\ell(\mathbf{u})\|_2 < \infty. \quad (6)$$

Throughout,  $\|\cdot\|_2$  denotes the Euclidean ( $l_2$ ) norm.

### C.2 POSITIONAL ATTENUATION

**Lemma 1** (Absolute positional embedding attenuates as  $m \rightarrow \infty$ ). *Let  $\mathbf{x}^{(0)} = m\mathbf{v} + \mathbf{p}$  with  $\|\mathbf{v}\|_2 = 1$ ,  $m > 0$ , and absolute positional embedding  $\mathbf{p} \in \mathbb{R}^d$ . Suppose  $\text{Norm} \in \{\text{LayerNorm}, \text{RMSNorm}\}$  and  $\mathbf{v}$  is non-degenerate for *LayerNorm* (i.e., its per-feature variance is nonzero; this holds for generic token embeddings). Then*

$$\|\text{Norm}(m\mathbf{v} + \mathbf{p}) - \text{Norm}(m\mathbf{v})\|_2 = \mathcal{O}\left(\frac{\|\mathbf{p}\|_2}{m}\right) \quad \text{as } m \rightarrow \infty \text{ (with } \mathbf{v}, \mathbf{p} \text{ fixed)}.$$

Hence the positional contribution shrinks linearly in  $1/m$ .

*Proof.* By scale invariance,  $\text{Norm}(m\mathbf{v} + \mathbf{p}) = \text{Norm}(\mathbf{v} + \varepsilon)$  with  $\varepsilon := \mathbf{p}/m$ , and  $\text{Norm}(m\mathbf{v}) = \text{Norm}(\mathbf{v})$ .

*RMSNorm.* With  $\|\mathbf{v}\| = 1$ ,

$$\frac{\mathbf{v} + \varepsilon}{\|\mathbf{v} + \varepsilon\|} = \mathbf{v} + (\mathbf{I}_d - \mathbf{v}\mathbf{v}^\top)\varepsilon + \mathcal{O}(\|\varepsilon\|^2),$$

540 hence  $\text{RMSN}(\mathbf{v} + \varepsilon) - \text{RMSN}(\mathbf{v}) = \sqrt{d}(\mathbf{I}_d - \mathbf{v}\mathbf{v}^\top)\varepsilon + \mathcal{O}(\|\varepsilon\|^2)$  and  
 541  $\|\text{RMSN}(m\mathbf{v} + \mathbf{p}) - \text{RMSN}(m\mathbf{v})\| \leq \sqrt{d} \|\mathbf{p}\| / m + \mathcal{O}(m^{-2})$ .  
 542

543 *LayerNorm.* Write  $\mathbf{a} := \mathbf{C}\mathbf{v} \neq \mathbf{0}$ ,  $\mathbf{u} := \mathbf{a} / \|\mathbf{a}\|$ . Then

$$544 \frac{\mathbf{a} + \mathbf{C}\varepsilon}{\|\mathbf{a} + \mathbf{C}\varepsilon\|} = \mathbf{u} + \frac{(\mathbf{I}_d - \mathbf{u}\mathbf{u}^\top)\mathbf{C}\varepsilon}{\|\mathbf{a}\|} + \mathcal{O}(\|\varepsilon\|^2),$$

545 so  $\|\text{LN}(m\mathbf{v} + \mathbf{p}) - \text{LN}(m\mathbf{v})\| = \sqrt{d} \frac{(\mathbf{I}_d - \mathbf{u}\mathbf{u}^\top)\mathbf{C}\mathbf{p}}{m\|\mathbf{C}\mathbf{v}\|} + \mathcal{O}(m^{-2})$ , which is  $\mathcal{O}(\|\mathbf{p}\| / m)$ .  $\square$   
 546  
 547

### 550 C.3 RESIDUAL STAGNATION

551 **Lemma 2** (Residual stagnation in a pre-norm block). *Let  $\mathbf{x}^{(\ell+1)} = \mathbf{x}^{(\ell)} + F_\ell(\text{Norm}(\mathbf{x}^{(\ell)}))$  with*  
 552  *$\mathbf{x}^{(\ell)} \neq \mathbf{0}$  and  $\text{Norm} \in \{\text{LN}, \text{RMSN}\}$ , and let*

$$553 B_\ell := \sup_{\mathbf{u} \in S} \|F_\ell(\mathbf{u})\|_2 < \infty.$$

554 Then

$$555 \frac{\|\mathbf{x}^{(\ell+1)} - \mathbf{x}^{(\ell)}\|_2}{\|\mathbf{x}^{(\ell)}\|_2} \leq \frac{B_\ell}{\|\mathbf{x}^{(\ell)}\|_2}, \quad \angle(\mathbf{x}^{(\ell)}, \mathbf{x}^{(\ell+1)}) \leq \arcsin\left(\frac{B_\ell}{\|\mathbf{x}^{(\ell)}\|_2}\right).$$

556 *Proof.* Since  $\text{Norm}(\mathbf{x}^{(\ell)}) \in S$ , we have  $\|\mathbf{x}^{(\ell+1)} - \mathbf{x}^{(\ell)}\|_2 = \|F_\ell(\text{Norm}(\mathbf{x}^{(\ell)}))\|_2 \leq B_\ell$ , giving  
 557 the first bound. Write  $\mathbf{x}^{(\ell+1)} = \mathbf{x}^{(\ell)} + \delta$ . The orthogonal component of  $\delta$  is at most  $\|\delta\|$ ; a short  
 558 calculation shows  $\sin \angle(\mathbf{x}^{(\ell)}, \mathbf{x}^{(\ell+1)}) \leq \|\delta\|_2 / \|\mathbf{x}^{(\ell)}\|_2 \leq B_\ell / \|\mathbf{x}^{(\ell)}\|_2$ , which implies the stated  
 559 angle bound.  $\square$   
 560

561 **Proposition 1** (Accumulated directional drift across  $L$  pre-norm blocks). *Let  $\mathbf{x}^{(0)} \neq \mathbf{0}$  and  $\mathbf{x}^{(\ell+1)} =$*   
 562  *$\mathbf{x}^{(\ell)} + F_\ell(\text{Norm}(\mathbf{x}^{(\ell)}))$  for  $\ell = 0, \dots, L-1$ . Let  $B_\ell := \sup_{\mathbf{u} \in S} \|F_\ell(\mathbf{u})\|_2 < \infty$ , and  $S_L :=$*   
 563  *$\sum_{j=0}^{L-1} B_j$ . Assume  $\|\mathbf{x}^{(0)}\|_2 > S_L$ , then*

$$564 \angle(\mathbf{x}^{(0)}, \mathbf{x}^{(L)}) \leq \frac{\pi}{2} \sum_{\ell=0}^{L-1} \frac{B_\ell}{\|\mathbf{x}^{(0)}\|_2 - \sum_{j<\ell} B_j} \leq \frac{\pi}{2} \frac{S_L}{\|\mathbf{x}^{(0)}\|_2 - S_L}.$$

565 *Proof.* Let  $\theta_\ell := \angle(\mathbf{x}^{(\ell)}, \mathbf{x}^{(\ell+1)})$ . By the recall above,  $\theta_\ell \leq \arcsin(B_\ell / \|\mathbf{x}^{(\ell)}\|) \leq \frac{\pi}{2} B_\ell / \|\mathbf{x}^{(\ell)}\|$ .  
 566 Also  $\|\mathbf{x}^{(\ell)}\| \geq \|\mathbf{x}^{(0)}\| - \sum_{j<\ell} B_j$  (each step can shrink the norm by at most  $B_\ell$ ). Summing angles  
 567 (spherical triangle inequality) gives the first display; since  $\|\mathbf{x}^{(0)}\| - \sum_{j<\ell} B_j \geq \|\mathbf{x}^{(0)}\| - S_L$ , each  
 568 fraction is  $\leq B_\ell / (\|\mathbf{x}^{(0)}\| - S_L)$ , yielding the last bound.  $\square$   
 569  
 570  
 571  
 572

## 573 D EMPIRICAL VALIDATION OF TOKEN GEOMETRY ANALYSIS

574 We empirically validate the two theoretical effects introduced in the previous sections. Effect I  
 575 describes the attenuation of positional information under large embedding magnitudes, while Effect II  
 576 concerns residual-update stagnation in pre-norm Transformer blocks. Our experiments directly probe  
 577 both behaviors on the base encoder, TI, and our proposed DTI.  
 578

579 **Effect I (Attenuation of positional information).** We evaluate whether increasing embedding  
 580 magnitude makes positional information unrecoverable after the first pre-norm normalization (LN).  
 581 To validate this, we train a 2-layer MLP classifier on the *frozen* base text encoder to predict a  
 582 token’s absolute position from the output of  $\text{LN}(\mathbf{e} + \mathbf{p})$ , where  $\mathbf{e}$  and  $\mathbf{p}$  each denotes token and  
 583 positional embeddings. On unmodified inputs (‘Normal’ in Figure 4), the classifier achieves 100%  
 584 accuracy, confirming that LN preserves positional information. We then scale the norm of a single  
 585 token embedding by a factor  $m \in \{0.5, 1, 2, 4, 8, 16\}$  before applying LN. Accuracy deteriorates  
 586 rapidly once  $m$  exceeds the natural scale of the encoder. Furthermore, we evaluated the classifier on  
 587 TI-trained and DTI-trained personalized embeddings. TI embeddings, which has excessively large  
 588 norms, collapse to near-zero positional accuracy, while DTI embeddings remain fully recoverable.  
 589  
 590  
 591  
 592  
 593

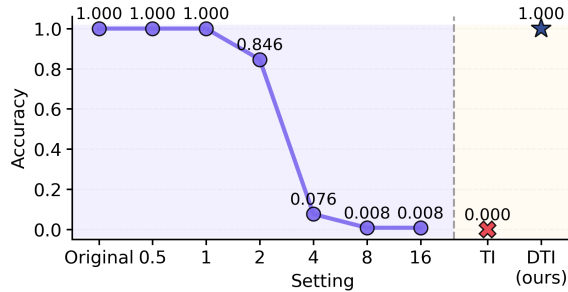


Figure 4: **Empirical validation.** Position prediction accuracy from LN outputs under varying embedding magnitudes, compared with trained TI and DTI embeddings.

This behavior directly corroborates Lemma 1: when  $m$  becomes large,  $\text{LN}(m\mathbf{v} + \mathbf{p})$  becomes dominated by  $m\mathbf{v}$ , rendering the positional component  $\mathbf{p}$  effectively invisible. DTI avoids this failure mode by constraining magnitudes to remain in-distribution.

**Effect II (Residual-update stagnation).** To test Lemma 2, we measure the internal angular change of hidden states within each pre-norm Transformer block. For each concept token, we compute the angle between the hidden state entering and exiting each block and then average across all layers. The average per-block angular change of TI embeddings was  $21.33^\circ$ , whereas the angular change of DTI embeddings was  $33.52^\circ$  ( $1.57\times$  larger).

These results support the theoretical prediction that excessively large norms suppress the effective residual direction in pre-norm blocks, causing the forward computation to behave nearly as an identity mapping. By keeping embedding norms within the training distribution, DTI prevents such stagnation and permits substantially larger and more meaningful updates throughout the encoder.

## E EXTENDED METHODS

### E.1 RSGD FOR TOKEN EMBEDDING OPTIMIZATION

We employ Riemannian Stochastic Gradient Descent (RSGD) to optimize the directional embedding  $\mathbf{v}$  while ensuring it remains on the unit hypersphere  $\mathbb{S}^{d-1}$ . The update step at iteration  $k$  involves projecting the Euclidean gradient onto the tangent space  $T_{\mathbf{v}_k}\mathbb{S}^{d-1}$  and retracting the updated vector back to the manifold:

$$\mathbf{g} = \mathbf{g}_{\text{euc}} - (\mathbf{v}_k^\top \mathbf{g}_{\text{euc}}) \mathbf{v}_k \in T_{\mathbf{v}_k}\mathbb{S}^{d-1}, \quad \mathbf{v}_{k+1} = \text{Retr}_{\mathbf{v}_k}(-\eta \mathbf{g}) = \frac{\mathbf{v}_k - \eta \mathbf{g}}{\|\mathbf{v}_k - \eta \mathbf{g}\|_2}. \quad (7)$$

Here,  $\mathbf{g}_{\text{euc}}$  is a Euclidean space gradient,  $\mathbf{g} \in T_{\mathbf{v}_k}\mathbb{S}^{d-1}$  is a tangent-space gradient, and  $\eta > 0$  is a learning rate. In practice, we scaled the gradient  $\mathbf{g}$  using its own norm similarly. This was inspired by Euclidean space optimizers (Hinton et al., 2012; Kingma & Ba, 2015; Loshchilov & Hutter, 2019), which normalizes the gradient based on moving average of squared gradients.

We observe that gradient magnitudes tend to increase as training progresses, which often leads to instability in the later stages. Although standard learning rate schedules can help mitigate this issue, the gradient dynamics vary considerably across different datasets and training settings, limiting the effectiveness of fixed schedules. To address this, we draw inspiration from adaptive optimization techniques in Euclidean space (Kingma & Ba, 2015; Duchi et al., 2011) and propose a simple yet effective gradient scaling scheme based on gradient norms:

$$\mathbf{g}'_k = \mathbf{g}_k / \|\mathbf{g}_k\|_2, \quad (8)$$

where  $\mathbf{g}_k$  is the gradient at iteration  $k$ . This normalization is equivalent to using an adaptive step size  $\eta / \|\mathbf{g}_k\|$  in a Euclidean update  $\mathbf{v}_{k+1} = \mathbf{v}_k - (\eta / \|\mathbf{g}_k\|_2) \mathbf{g}_k$ ; the update direction is still  $\mathbf{g}_k$ , but the length is fixed to  $\eta$ , preventing very large gradients from causing excessively large parameter updates. Note that a similar technique was previously explored in the context of Riemannian optimization (Cho & Lee, 2017).

## E.2 WHY vMF OVER OTHER DISTRIBUTIONS?

We chose the von Mises-Fisher (vMF) distribution as it is ideally suited for modeling the directional characteristics of token embeddings we identified in Section 2. Our central hypothesis is that the token embedding vocabulary can be modeled as a **mixture of vMF distributions**, where each component corresponds to a distinct semantic cluster (e.g., one for animals, another for objects). The vMF distribution is the suitable building block for this model for three key reasons:

- **It’s a natural fit.** The vMF is the natural analog to the Gaussian distribution on a hypersphere, making it a principled and standard choice for modeling directional data clusters.
- **It’s computationally efficient.** The vMF’s mathematical form is exceptionally convenient for optimization. In our MAP formulation, the gradient of the log-prior is a *constant-direction vector* ( $-\kappa\mu$ ), which provides a stable and efficient semantic pull without requiring complex calculations. This simplicity makes it more suitable for high-dimensional embeddings in large-scale models than alternatives like the Kent and Bingham distributions.
- **It’s interpretable and controllable.** The parameters are easy to understand. The mean direction  $\mu$  serves as a *semantic anchor* to prevent the learned token from drifting away from related concepts, while the concentration  $\kappa$  allows us to control the strength of this regularization.

These factors collectively make the vMF distribution a superior choice for our application, providing the necessary regularization in a way that is both mathematically principled and computationally tractable.

## F EXTENDED EXPERIMENTS

---

### Algorithm 1 Directional Textual Inversion (DTI)

---

```

1: Inputs: Model  $\epsilon_\theta$ , text encoder  $c(\cdot)$ , init token  $e_{\text{init}}$ , magnitude  $m^*$ ,  $\kappa$ , iterations  $K$ , learning rate  $\eta$ 
2:  $v_0 \leftarrow e_{\text{init}} / \|e_{\text{init}}\|_2$ 
3:  $\mu \leftarrow e_{\text{init}} / \|e_{\text{init}}\|_2$ 
4: for  $k = 0$  to  $K - 1$  do
5:   Sample minibatch  $(z, t, \epsilon)$ 
6:    $g_{\text{data}} \leftarrow \nabla_v \mathcal{L}_{\text{data}}(m^* v_k)$ 
7:    $g_{\text{euc}} \leftarrow g_{\text{data}} - \kappa \mu$  (add prior gradient)
8:    $g \leftarrow g_{\text{euc}} - (g_{\text{euc}}^\top v_k) v_k$  (tangent projection)
9:    $g' \leftarrow g / \|g\|_2$  (gradient scaling)
10:   $v_{k+1} \leftarrow \frac{v_k - \eta g'}{\|v_k - \eta g'\|_2}$  (retraction to  $\mathcal{S}^{d-1}$ )
11: end for
12: return  $e^* = m^* v_K$ 

```

---

### F.1 IMPLEMENTATION DETAILS

Following the protocol of recent studies, we primarily conducted experiments using Stable Diffusion XL (SDXL). To demonstrate broader applicability to different models, we also conducted experiments with very recent model, SANA 1.5 (Xie et al., 2025), where the results can be found in Table 1.

For a fair comparison, we adopted most of the hyperparameter settings from the Textual Inversion (TI) implementation provided by the HuggingFace `diffusers` library. Specifically, we used a training batch size of 4, and enabled mixed-precision training with the `bfloat16` (bf16) format. We set the learning rate commonly-used  $5e-3$ . All experiments were run with a fixed random seed of 42, and the maximum number of training steps was set to 500. For output generation, we used the `DDIMScheduler` with 50 inference steps for SDXL and 20 steps with `FlowMatchEulerDiscreteScheduler` for SANA.

**Hyperparameters.** There can be various approaches to selecting the concentration parameter  $\kappa$ . We performed a grid search and found that values in the range of  $5e-5$  to  $2e-4$  works well. Therefore, we did not conduct an extensive search for an optimal decimal value. Throughout the experiments, we simply fixed value to  $1e-4$ , which generalizes well to experiments with different settings. Examples illustrating the effects of different  $\kappa$  settings are provided in Table 2.

**Baselines.** Throughout this paper, we compare our method with two baseline approaches: Textual Inversion (TI) (Gal et al., 2023a) and CrossInit (Pang et al., 2024a). Since the official CrossInit implementation is based on Stable Diffusion v2.1 with hyperparameters tailored to that version, we reconfigure it to operate on SDXL by aligning its training setup with that of TI. Specifically, we adopt the same hyperparameters as used for TI, and we set the regularization weight for CrossInit to  $1e-5$ , as specified in the original paper.

## F.2 ON THE CHOICE OF PRIOR

For all of our experiments in the main section, we used the initial tokens as prior from the DreamBooth dataset as is. However, we would like to note that since our DTI can leverage the prior, searching for better priors can lead to better results. This demonstrates the effectiveness of the prior.

To test this, we experimented with having a VLM recommend initial tokens. More specifically, we provided reference images to the VLM and asked it to recommend 1-2 words that best describe them. For the experiments, we used Qwen-VL 2.5 (Bai et al., 2025) as the VLM. The results are shown in Table 4.

The results indicate that changing the prior affects performance, although the overall effect is modest. For both TI and our DTI, Qwen-VL initialization tends to increase subject similarity, accompanied by a slight decrease in text fidelity. Practitioners may leverage VLMs or manually craft priors with targeted terms to emphasize desired attributes. Overall, these findings demonstrate the flexibility and effectiveness of leveraging priors.

Table 4: Results with VLM-recommended priors. We compare Qwen-VL recommended initial tokens with DreamBooth initial tokens as priors for DTI.

Method	Initialization	SDXL		SANA	
		Image	Text	Image	Text
TI	DreamBooth init	0.561	0.292	0.480	0.621
	Qwen-VL init	0.583	0.273	0.501	0.619
DTI (ours)	DreamBooth init	0.450	0.522	0.479	0.744
	Qwen-VL init	0.520	0.391	0.504	0.697

## F.3 COMPARISON WITH OTHER BASELINES

We expand our comparative analysis to include additional baselines: P+ (Voyunov et al., 2023), NeTI (Alaluf et al., 2023), and CoRe (Wu et al., 2025). We run these experiments mainly on SD1.5 and SD2.1-base as these baseline papers work on those versions. Adhering to the evaluation protocol of the main paper, we measure subject similarity using DINOv2 similarity and prompt fidelity with the CLIP-variant, SigLIP. The results demonstrate that across both architectures, DTI consistently achieves the most favorable balance between these metrics compared to all baselines. Qualitative comparison can be found in Figure 9.

**DTI as a drop-in replacement for TI.** Although DTI is primarily designed for embedding-only personalization, it also functions effectively as a drop-in replacement within fine-tuning pipelines. Recent work on Direct Consistency Optimization (DCO) Lee et al. (2024) typically performs joint optimization of a concept token using standard Textual Inversion (TI). To assess the impact of substituting this TI component with DTI, we conducted joint training for 250 steps using a LoRA module with rank 4.

As shown in Table 6, the conventional TI-based joint training exhibits limited text–alignment (0.456), whereas replacing TI with DTI substantially improves alignment to 0.635. This quantitative improve-

Table 5: **Results on SD1.5 and SD2.1-base.** We compare the baselines that improve TI on different versions of Stable Diffusion. DTI achieves the best balance between subject similarity and text fidelity compared to other baselines.

Method	SD1.5		SD2.1-base	
	Image	Text	Image	Text
P+ (Voyunov et al., 2023)	0.273	<b>0.719</b>	0.238	<b>0.663</b>
NeTI (Alaluf et al., 2023)	0.408	0.579	0.565	0.517
CoRe (Wu et al., 2025)	0.340	0.661	0.357	0.654
DTI (ours)	<b>0.418</b>	0.554	<u>0.469</u>	0.568

ment is further reflected in Figure 5, where our method accurately incorporates textual attributes (e.g., red backpack), while DCO with TI fails to do so.

Table 6: **DCO experiments.** Comparison of DCO with standard Textual Inversion (TI) versus DCO initialized with our DTI. Integrating our method significantly improves text alignment.

Method	Image	Text
DCO	<b>0.605</b>	0.456
DCO + DTI (ours)	0.568	<b>0.635</b>

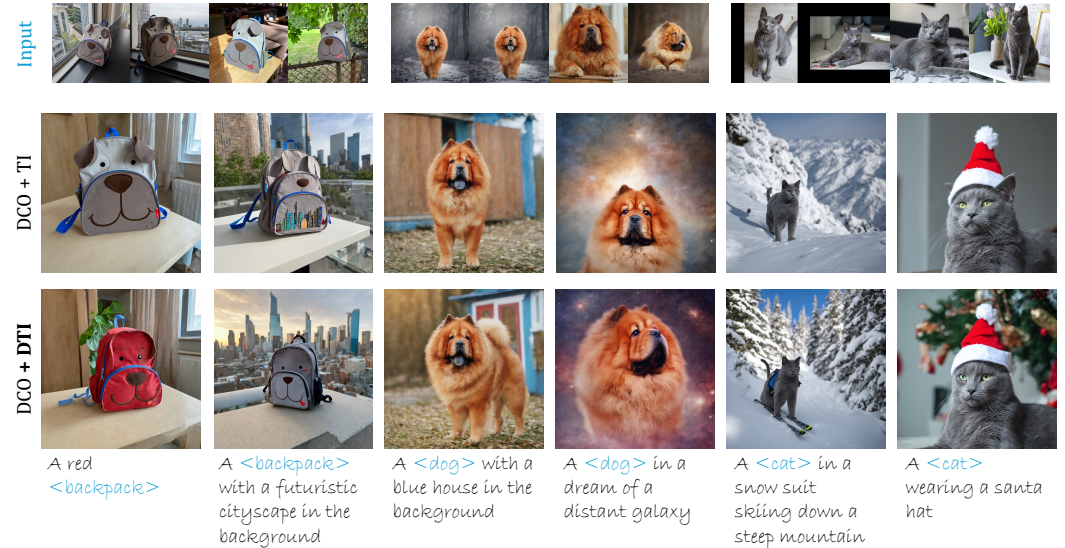


Figure 5: **Qualitative results with DCO (Lee et al., 2024).** We provide the comparison of the output images of DCO + TI and DCO + DTI. The results suggest our DTI’s superiority in text fidelity while preserving strong image similarity.

#### F.4 ABLATION STUDY

**Effect of Riemannian optimization.** Our DTI framework employs Riemannian optimization to ensure embeddings lie on the spherical manifold  $\mathbb{S}^{n-1}$ . An alternative is to simply re-scale embeddings after each Euclidean optimization step to achieve this constraint. However, Table 2 (first row) shows this latter Euclidean-based approach with re-scaling achieves suboptimal results, highlighting the benefit of direct Riemannian optimization.

**Effect of magnitude ( $m$ ).** We investigated the impact of the fixed embedding magnitude,  $m$ , on personalization performance. Our DTI framework, by default, sets  $m$  to the average norm observed in the pre-trained CLIP token vocabulary. We compared this “mean” strategy under the Riemannian optimization setting with  $\kappa = 1e - 4$ :

Table 7: We surveyed real-world user preferences regarding subject fidelity and image-text alignment. DTI ranks the top in both metrics, confirming its practical benefits.

	TI	CrossInit	DTI (ours)
Image fidelity	13.78	42.87	<b>43.45</b>
Text alignment	10.83	22.40	<b>66.77</b>

- Setting  $m$  to the minimum vocabulary norm (“min”).
- Setting  $m$  to the mean vocabulary norm (“mean”).
- Setting  $m$  to a large, out-of-distribution (OOD) value of 5.0.

As shown in Table 2:

- The “mean” strategy achieves the highest subject similarity and strong text fidelity.
- The “min” strategy results in significantly poorer performance in both metrics.
- Using an OOD magnitude of 5.0 also leads to a degradation in both metrics.

These results validate our choice of fixing the magnitude to an in-distribution scale, specifically the average vocabulary norm, as it provides a strong balance of subject similarity and text alignment. Both excessively small (“min”) and out-of-distribution large (“OOD”) magnitudes are detrimental.

**Effect of concentration parameter ( $\kappa$ ).** The concentration parameter  $\kappa$  of the von Mises-Fisher (vMF) prior controls the strength of the directional regularization. We analyzed its effect by varying  $\kappa$  while using Riemannian optimization and the “mean” embedding magnitude. We tested  $\kappa = 0.0$  (no prior),  $\kappa = 1e - 4$  (DTI default), and  $\kappa = 5e - 4$ .

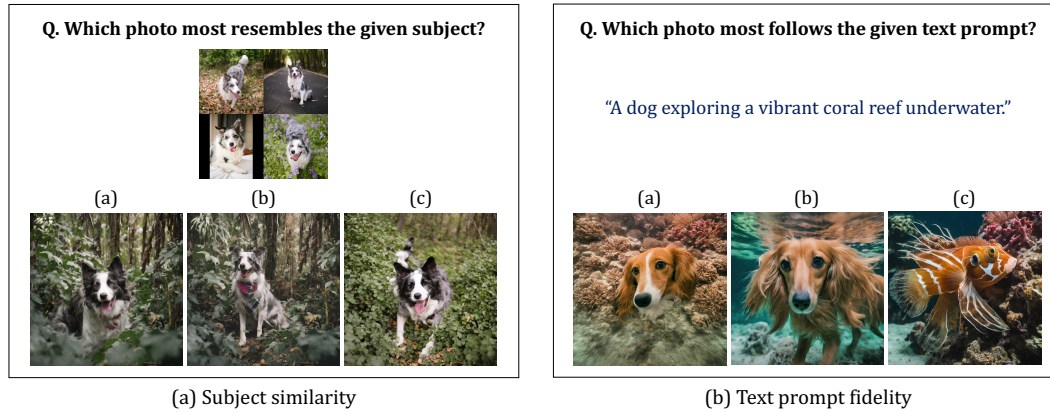
The results in Table 2 indicate:

- With  $\kappa = 1e - 4$ , we observe the best balance between subject similarity and text fidelity.
- Setting  $\kappa = 0.0$ , which removes the directional prior, leads to lower scores in text fidelity, which validates our method’s priority in model’s enhancing semantic understanding.
- Increasing the regularization strength with  $\kappa = 5e - 4$  yields the highest text fidelity among the tested values but at the cost of reduced subject similarity.

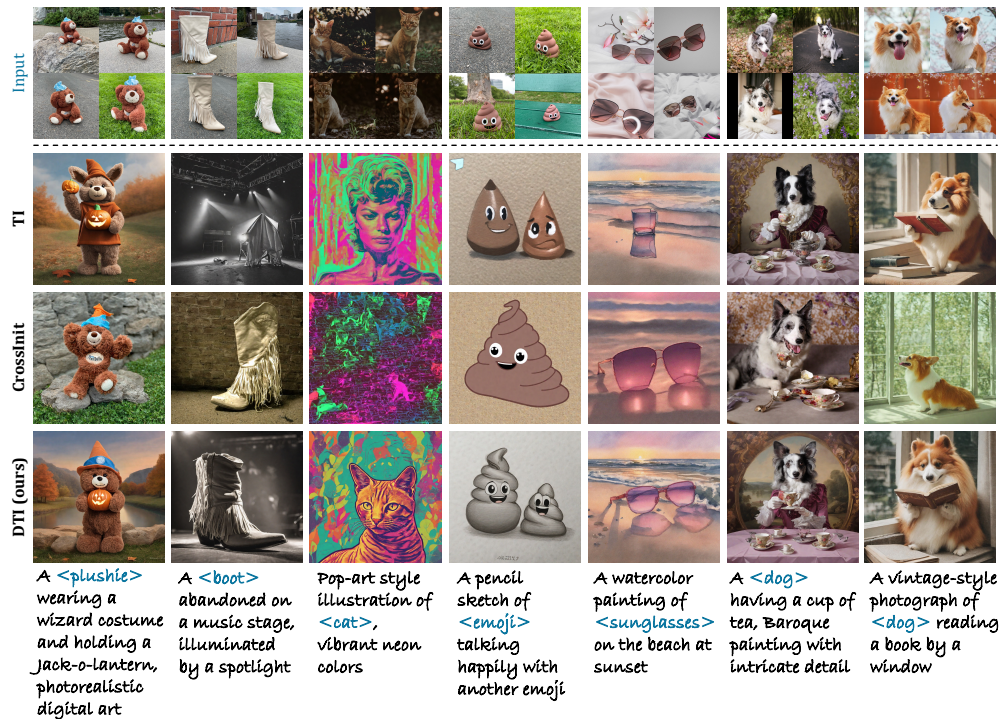
Overall, our default choice of  $\kappa = 1e - 4$  provides a better balance between maintaining subject similarity and ensuring text fidelity. Note that  $\kappa = 1e - 4$  may not be strictly optimal in decimals across all criteria but works reasonably well by providing robust overall performance.

## F.5 USER STUDY

To evaluate real-world user preferences for image generation quality, we conducted a comprehensive user study involving 100 participants recruited through Amazon Mechanical Turk. Each participant completed a survey consisting of 20 questions, evenly divided into two critical evaluation criteria: subject similarity and text prompt fidelity. For each question, participants were presented with three distinct image options, generated by: Textual Inversion (Gal et al., 2023a), CrossInit (Pang et al., 2024a), and our proposed Directional Textual Inversion (DTI). The order of these three choices was randomized for each question, using a fixed random seed to ensure consistent shuffling across all participants. Sample questions can be found in Figure 6. We collected a total of 96 valid responses, with 4 submissions being excluded due to invalid patterns such as selecting the same answer for all questions. The results, as detailed in Table 7, demonstrate that our Directional Textual Inversion (DTI) consistently outperforms both Textual Inversion and CrossInit across both evaluation metrics: image subject similarity and text prompt fidelity. These findings confirm the superior performance of our proposed method in generating images that more accurately align with user expectations regarding both visual content and textual descriptions.



878 **Figure 6: User study design.** We conducted a user study with 100 participants recruited via Amazon  
 879 Mechanical Turk to evaluate 20 questions. The evaluation focused on two key aspects: subject  
 880 similarity (10 questions) and text prompt fidelity (10 questions). To ensure fair comparison, the  
 881 random seed was fixed and option order was shuffled.



906 **Figure 7: Qualitative results on SDXL.** We compare DTI with previous methods across diverse  
 907 subjects and textual prompts, spanning simple descriptions to complex variations in attributes,  
 908 backgrounds, and styles (same random seeds). All results in this figure are generated with SDXL  
 909 (SANA in Appendix Figure 8).

## 912 F.6 QUALITATIVE RESULTS

914 Figure 7 illustrates qualitative comparisons across various prompts. DTI consistently generates  
 915 images that more accurately reflect the content of the captions, while effectively preserving subject  
 916 consistency. For instance, for 'Pop-art style illustration of <cat>', TI omits the cat while DTI renders  
 917 the cat in the specified style. Similarly, in the second column, TI and CrossInit fail to incorporate all  
 elements of the prompt, disregarding either the subject or details such as 'music stage' and 'spotlight'.

918 In contrast, DTI integrates both the subject and these details, producing a more complete output.  
919 Collectively, these examples highlight DTI’s superior compositional fidelity and subject preservation,  
920 showing its powerfulness that consistently satisfies all prompt constraints. This attributes to DTI’s  
921 stable optimization within the directional space, which facilitates improved integration of multiple  
922 prompt components. DTI’s ability to maintain subject fidelity and adhere to textual intent establishes  
923 it as a robust choice for a wide range of text-to-image generation tasks. Additional qualitative results  
924 including those of SANA can be found in Appendix F.6.

925 We present additional qualitative comparisons with TI-based approaches (Gal et al., 2023a; Pang et al.,  
926 2024a) in Figure 7 (SDXL) and 8 (SANA). The results illustrate that our proposed DTI consistently  
927 generates outputs that accurately align with the provided text prompts, even in challenging cases  
928 where the baseline methods fail to do so.

929 Our DTI serves as a drop-in replacement for TI, enhancing the model’s performance when combined  
930 with LoRA. The qualitative results in Figure 10 demonstrate that DTI consistently generates outputs  
931 that both precisely follow the text prompt and accurately capture the subject’s details.

## 933 F.7 MORE RESULTS ON APPLICATIONS

934  
935 **Stylization.** We explore the combination of personalized subject embeddings and style embeddings.  
936 Our method, DTI, consistently generates images that accurately reflect both the personalized subject  
937 and the specified style. In contrast, TI frequently fails in this task, either by omitting the subject  
938 altogether (top row) or by inadequately capturing the intended style or subject details (bottom row) of  
939 Figure 11.

940 **My object in my style.** We also compare our results in simultaneous generation of personalized  
941 subject and style. The results demonstrated in Figure 12 shows that DTI successfully generates  
942 outputs that are faithful to both subject and style, while TI fails to.

943 **Face personalization.** To evaluate and showcase the capability of our DTI method in face personal-  
944 ization, we conducted experiments using randomly selected faces from the FFHQ dataset (Karras  
945 et al., 2019) as well as faces generated by DALL·E (Ramesh et al., 2021).

946 Since CrossInit specifically focuses on facial personalization, we compare TI, CrossInit and our DTI  
947 on this task. Given that CrossInit does not explicitly provide hyperparameters (including learning rate)  
948 tailored for SDXL, we performed a grid search across various learning rates. Our empirical results  
949 indicated that the learning rate used by TI yielded reasonable performance for CrossInit as well.  
950 Figure 14 illustrates a comparison between the three methods, demonstrating that all methods perform  
951 effectively for facial personalization. Nevertheless, as the complexity of text prompts increases  
952 (rows depicted in the left columns), the baseline methods struggle to accurately reflect all described  
953 components of the prompts. In contrast, our DTI method consistently captures the critical components  
954 precisely, demonstrating superior performance in achieving enhanced textual fidelity.

955 **More results on interpolation.** More results on interpolating the personalized subjects can be found  
956 in Fig. 13, demonstrating our creative application capability.

## 958 G ADDITIONAL EXPERIMENTS

959  
960 We present additional experimental results in Figures 15, 16, 17, and 18. Specifically, Figure 15  
961 compares TI using SLERP against LERP, justifying our choice of the latter. In Figure 16, we present  
962 an ablation study on magnitude settings. While our DTI uses the mean value of the entire vocabulary  
963 as the default, we further investigate initializing with the specific category describing the subject  
964 (e.g., cat). We demonstrate that minor variations in magnitude do not significantly alter the outcome.  
965 Figure 17 evaluates our DTI in multi-concept scenarios, illustrating both successful outcomes and  
966 limitations. Finally, we analyze specific failure cases of our method in Figure 18.

## 968 H SOCIETAL IMPACTS

969  
970 The rapid advancement of text-to-image diffusion models, especially in the domain of personalization  
971 techniques, raises important societal considerations. In particular, the ease of generating highly

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

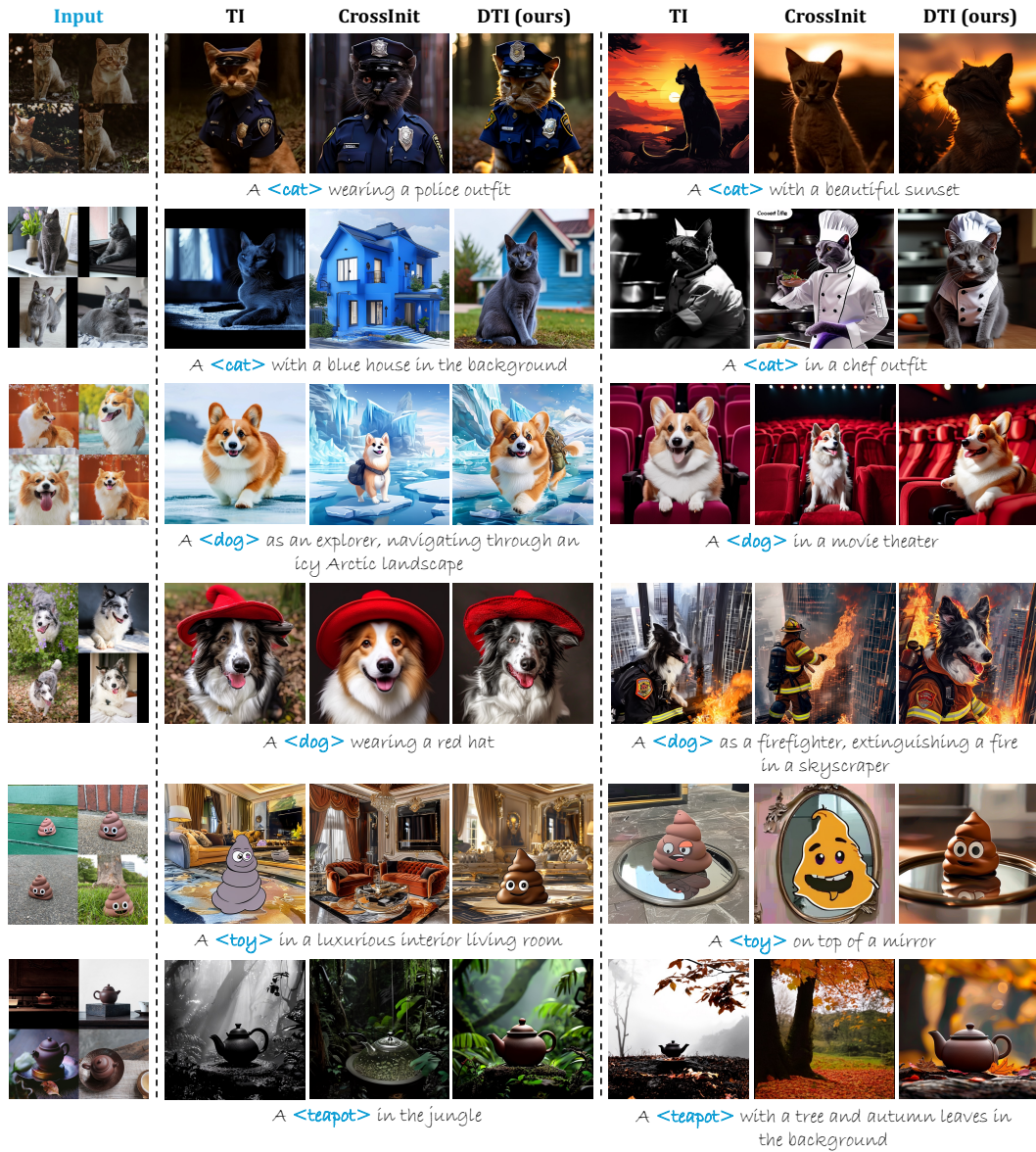
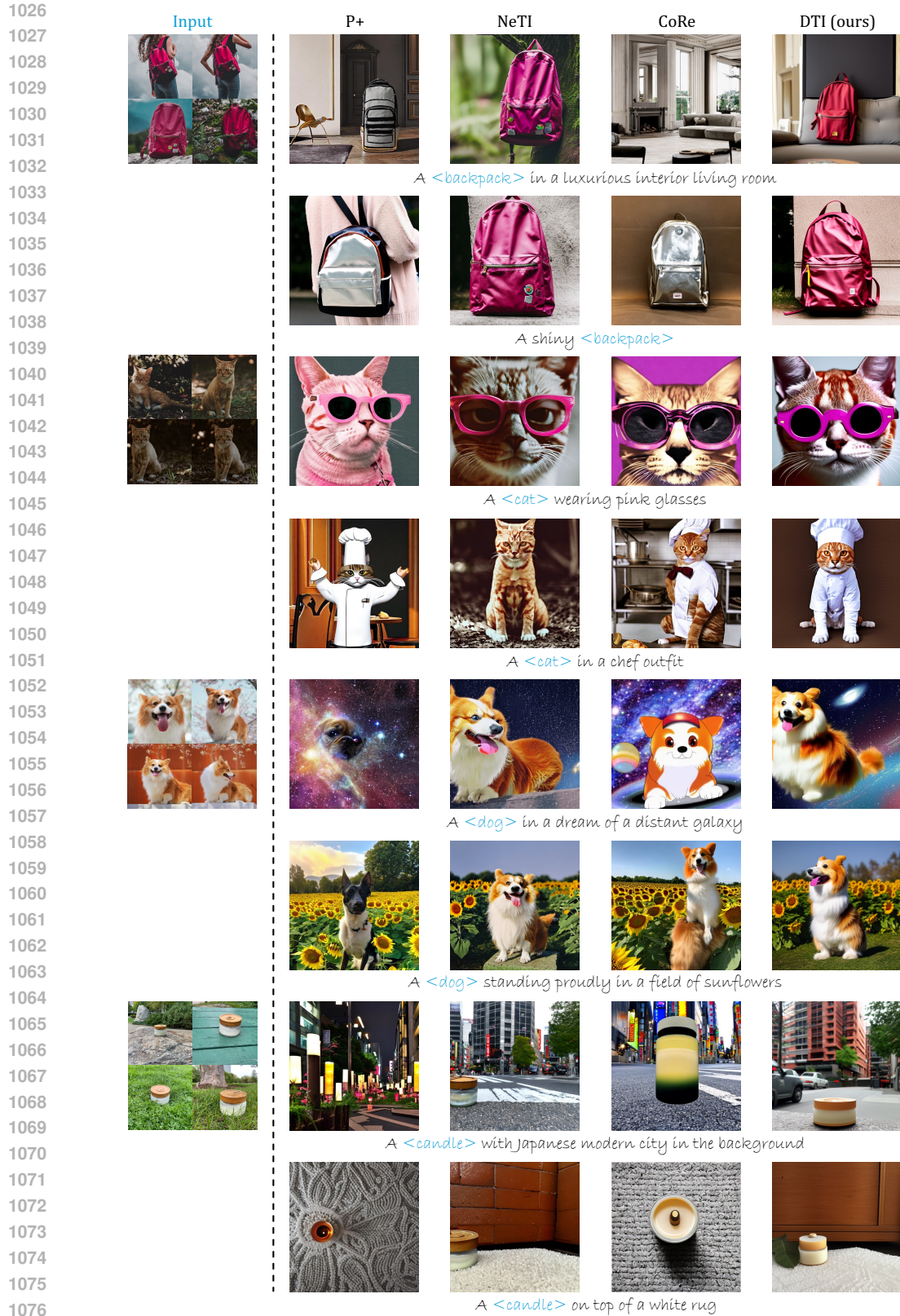


Figure 8: **Qualitative results with SANA1.5-1.6B.** Here, we provide more qualitative comparison with TI and CrossNit on SANA. Our DTI consistently generates results that precisely reflect the user text prompts, maintaining the subject similarity at the same time.



1078 **Figure 9: Comparison with additional baselines.** We provide qualitative comparisons against  
 1079 additional TI-enhancing methods—P+, NeTI, and CoRe. Because these baselines are built on SD2.1-  
 base, we apply DTI using the same pre-trained backbone to ensure fairness. The results demonstrate  
 that DTI attains higher text fidelity while maintaining subject similarity.

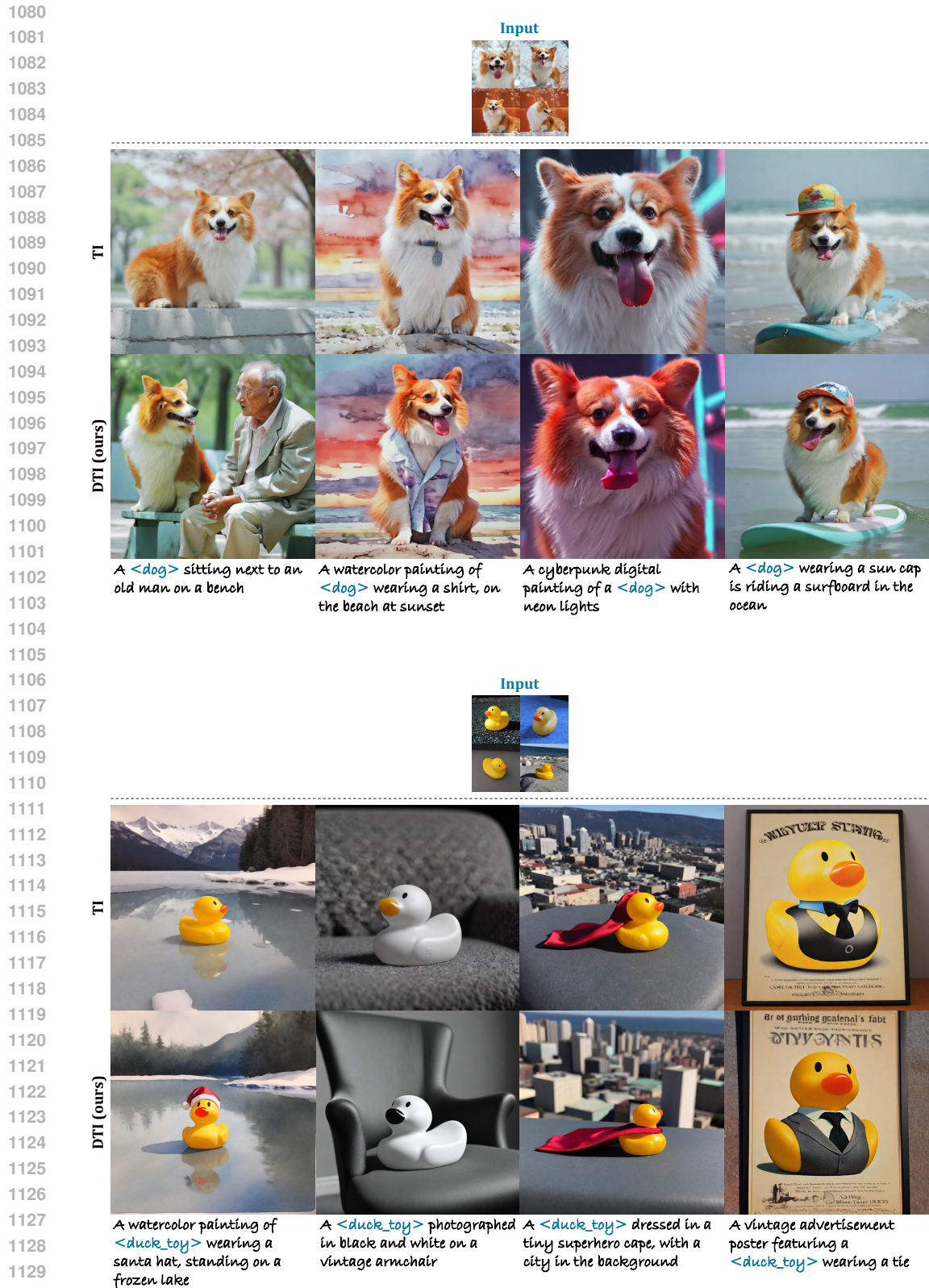


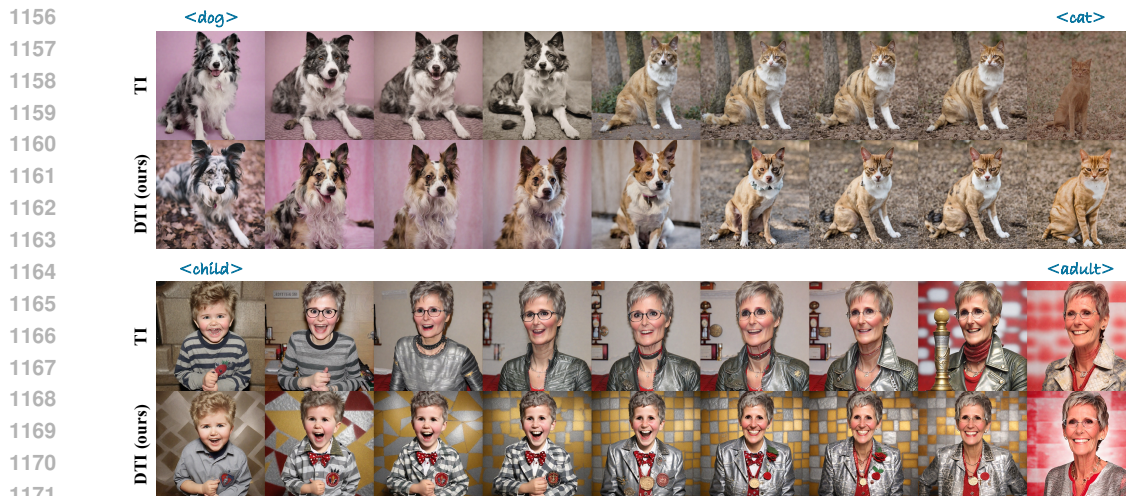
Figure 10: **Qualitative results with TI/DTI with LoRA on SDXL.** We have performed qualitative comparison of applying TI and DTI on model fine-tuning methods using LoRA (rank 32). DTI consistently improves the text prompt fidelity compared to TI.



1152 Figure 11: **Stylization.** Qualitative comparison of  
1153 personalization with diverse style inputs.



1152 Figure 12: **My subject in my style.** Qualitative  
1153 comparison of subject-style mixing within the  
1154 same prompt.



1172 Figure 13: **Interpolation.** We compare images generated by a TI and our DTI. Two per-  
1173 sonalized subjects are interpolated, including interpolation between inanimate and animate  
1174 subjects, live subjects, and human faces. Images are generated with interpolation ratio  
1175 [0.0, 0.35, 0.40, 0.45, 0.50, 0.55, 0.60, 0.65, 1.0] for better visualization. Our DTI offers smooth  
1176 interpolation between concepts, expanding the personalization in more creative axis.

1177  
1178  
1179 specific and detailed images can raise concerns related to copyright infringement, as personalized  
1180 generative models may inadvertently or intentionally reproduce objects protected by intellectual  
1181 property laws. Therefore, we note that it is important for users and distributors of the model to  
1182 develop comprehensive awareness and implement guidelines addressing copyright boundaries, fair  
1183 use, and ethical content generation. Moreover, we note that, since our method does not modify the  
1184 underlying parameters of the generative model but solely adjusts the token embeddings that capture  
1185 personalized concepts, the quality of generated images inherently depends on the capabilities of the  
1186 underlying text-to-image model.

1187

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

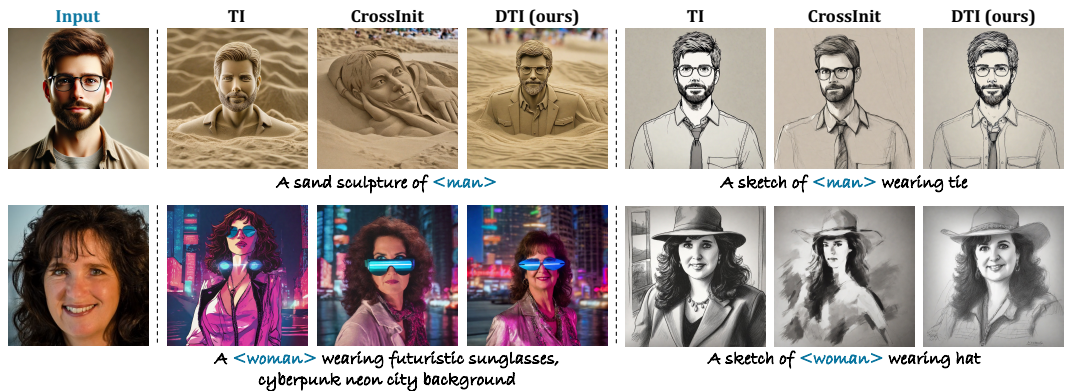


Figure 14: **Comparison of face personalization methods.** We compare our method and Textual Inversion (TI) against CrossInit, which specifically targets face personalization. To prevent bias from celebrity faces, we evaluate personalization using two alternative sources: images generated by DALL·E (Ramesh et al., 2021) (top row) and randomly selected images from the FFHQ (Karras et al., 2019) (bottom row).

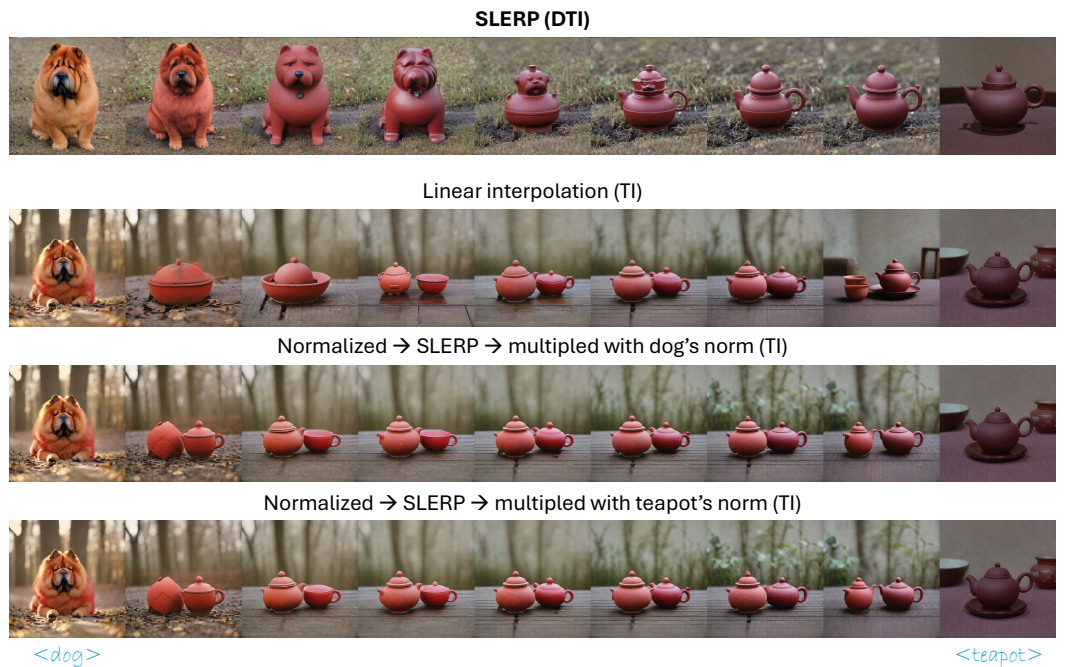
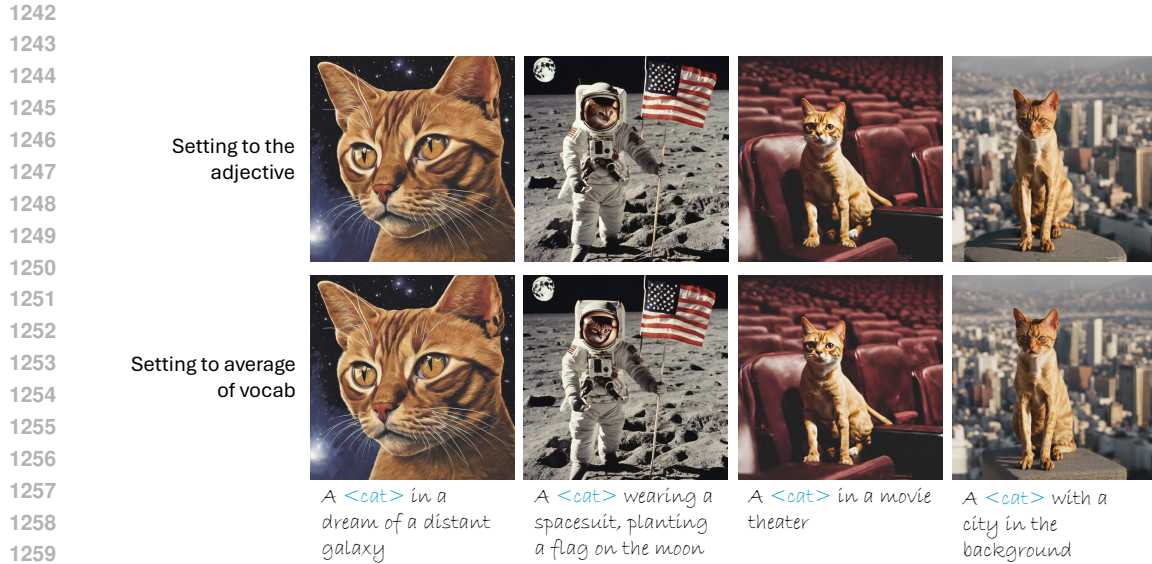
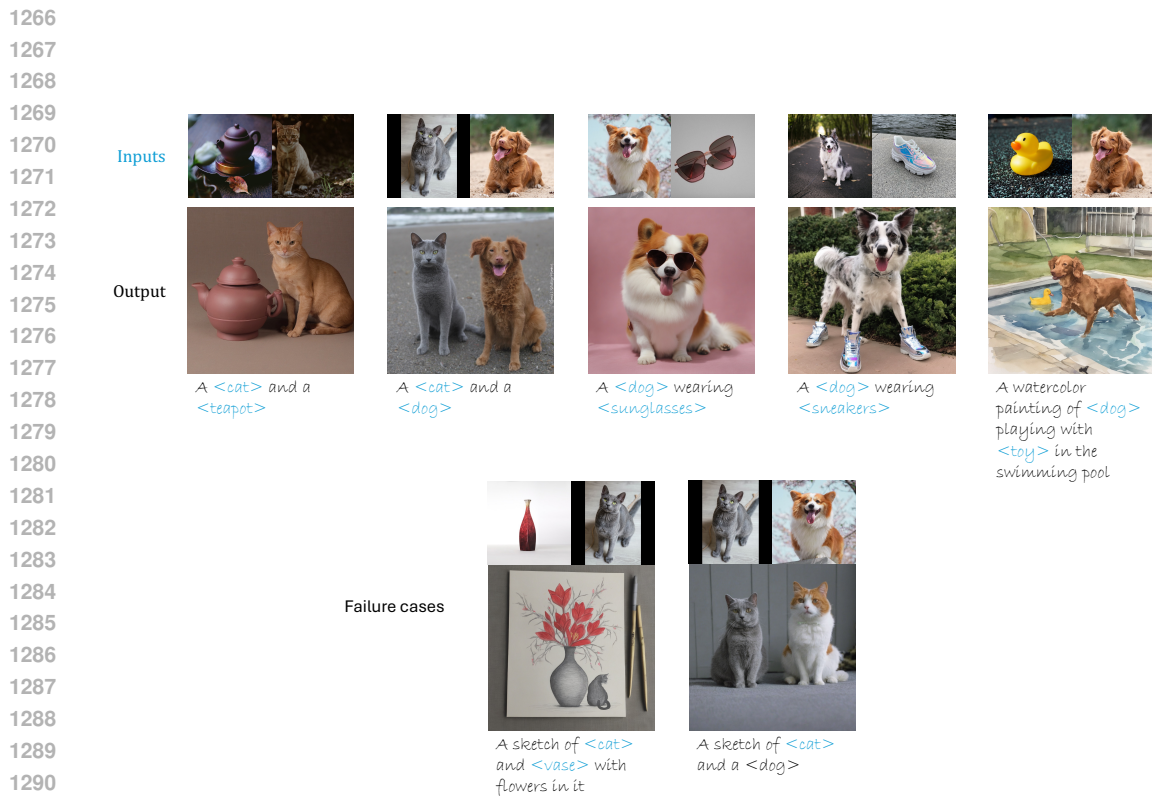


Figure 15: **Interpolation options for TI.** We compare several interpolation options for TI, including linear interpolation, SLERP with normalization and adjusted norms. While these approaches exhibit minor differences in behavior, none produce smooth transitions between concepts. In contrast, our DTI with SLERP achieves noticeably smoother and more consistent interpolations.



1261 **Figure 16: Ablation on magnitude settings.** For consistency and ease of use, all magnitudes in this  
1262 paper are set to the average value computed over the model’s vocabulary (see ablation in Table 3). To  
1263 evaluate the effect of using concept-specific magnitudes (e.g., the magnitude of ‘cat’ for the concept  
1264  $\langle \text{cat} \rangle$ ), we provide ablation results under different magnitude settings. The results show that small  
1265 deviations from the default magnitude do not lead to noticeable differences in output quality.  
1266

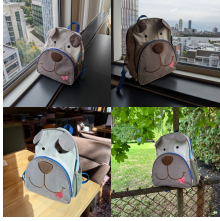


1292 **Figure 17: Multi-concept experiments.** We further evaluate DTI by combining multiple learned  
1293 concepts within a single prompt. The results demonstrate that DTI can successfully integrate multiple  
1294 concepts, while the second column shows failure cases exhibiting attribute binding issues.  
1295

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

**When subject requires high visual detail**

Input



A `<backpack>` with a mountain in the background



A `<backpack>` on top of a wooden floor



A `<backpack>` with a futuristic cityscape in the background



A `<bear_plushie>` besides a rushing waterfall



A `<bear_plushie>` in a field of wildflowers



A `<bear_plushie>` in a luxurious living room

**When prompt is difficult to depict**



A `<subject>` floating in an ocean of milk



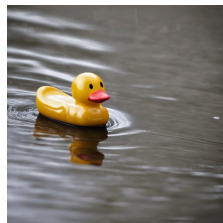
A cube shaped `<subject>`



**When prompt requires change of attributes**



A wet `<subject>`



A purple `<subject>`

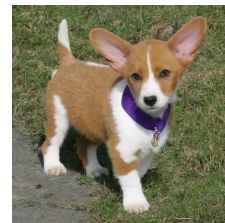


Figure 18: **Failure cases.** We present examples of three representative failure modes: (1) subjects that require high visual detail, (2) prompts that are vague or difficult to faithfully depict, and (3) prompts that involve attribute modification (e.g., color changes).