

000 SUPPLEMENTARY FILE  
 001  
 002 FLASHEDIT: DECOUPLING SPEED, STRUCTURE,  
 003 AND SEMANTICS FOR PRECISE IMAGE EDITING  
 004  
 005

006 **Anonymous authors**

007 Paper under double-blind review  
 008  
 009  
 010  
 011

012 1 A COMPREHENSIVE OVERVIEW OF DIFFUSION MODELS  
 013

014 Positioned at the forefront of generative modeling, Diffusion Probabilistic Models (DMs) (Ho et al.  
 015 (2020); Nichol & Dhariwal (2021); Kingma et al. (2021); Karras et al. (2022); Song et al. (2021);  
 016 Chen et al. (2023); Salimans & Ho (2022)) represent a significant paradigm shift, celebrated for  
 017 their exceptional capacity to generate high-fidelity outputs. Unlike Generative Adversarial Net-  
 018 works (GANs), which often suffer from training instability, or Variational Autoencoders (VAEs),  
 019 which can produce blurry results, DMs offer stable training and unparalleled sample quality. Their  
 020 influence spans numerous applications, from image and video creation to text-to-image synthesis  
 021 and audio generation. By providing a powerful, principled methodology for learning complex data  
 022 distributions, these models have become a fundamental pillar of contemporary artificial intelligence.  
 023 More recently, the fusion of diffusion principles with Transformer architectures has given rise to Dif-  
 024 fusion Transformers (DiTs) (Yang et al. (2025); Peebles & Xie (2023); Labs (2024); Hatamizadeh  
 025 et al. (2024); Lin et al. (2022); Kim et al. (2024)), which demonstrate remarkable scaling properties.  
 026 This section offers a thorough examination of their core concepts, mathematical underpinnings, and  
 027 architectural variations.

028 1.1 CORE PRINCIPLES OF DIFFUSION MODELS  
 029

030 Conceptually rooted in non-equilibrium thermodynamics and statistical mechanics, DMs are a class  
 031 of latent variable models. They operate through a carefully defined two-stage process that mirrors  
 032 the physical phenomenon of diffusion. The first stage, the **forward process**, systematically de-  
 033 grades structured data (a low-entropy state) by incrementally adding Gaussian noise over a series of  
 034 timesteps until only pure, unstructured noise remains (a high-entropy state). The second stage, the  
 035 **reverse process**, learns to invert this procedure. Here, a neural network is trained to methodically  
 036 remove the noise, starting from a random sample and progressively refining it to reconstruct a sam-  
 037 ple from the original data distribution, effectively decreasing the system’s entropy. This approach  
 038 transforms the intractable problem of modeling a complex data manifold into a more manageable  
 039 sequence of denoising steps. The power of DMs originates from this iterative refinement, which  
 040 enables the generation of highly detailed and coherent samples.

041 1.2 THE FORWARD NOISING PROCESS  
 042

043 The forward process Ho et al. (2020) is a fixed (non-learned) Markov chain that gradually perturbs  
 044 an initial data sample  $x_0 \sim q(x_0)$  over  $T$  discrete timesteps. At each step  $t$ , a small amount of  
 045 Gaussian noise is added according to a variance schedule  $\beta_t$ :

$$046 q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I). \quad (1)$$

047 The sequence  $\{\beta_t\}_{t=1}^T$  is a critical hyperparameter, typically chosen such that  $\beta_1 < \beta_2 < \dots <$   
 048  $\beta_T$ . Common schedules include linear, quadratic, or cosine schedules, which control the rate of  
 049 information destruction. As  $t \rightarrow T$ , the cumulative effect of the noise addition ensures that  $x_T$   
 050 converges to an isotropic Gaussian distribution,  $x_T \sim \mathcal{N}(0, I)$ .

051 A key property of this process is that we can sample  $x_t$  directly from  $x_0$  at any timestep  $t$  without  
 052 iteration. Let  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ . Using the reparameterization trick, we can write  $x_t$   
 053 as a function of  $x_0$  and a random noise variable  $\epsilon \sim \mathcal{N}(0, I)$ :

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon. \quad (2)$$

This gives the closed-form sampling distribution:

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I). \quad (3)$$

Here,  $\sqrt{\bar{\alpha}_t}$  represents the signal rate, while  $\sqrt{1 - \bar{\alpha}_t}$  is the noise rate. This property is instrumental for efficient model training.

### 1.3 THE LEARNED REVERSE DENOISING PROCESS

The reverse process is the generative core of the model, tasked with inverting the forward noising process. It begins with a sample from the prior distribution,  $x_T \sim \mathcal{N}(0, I)$ , and aims to learn the true posterior distribution  $q(x_{t-1}|x_t)$  to iteratively denoise the sample back to  $x_0$ . This reverse process is modeled as a Markov chain with learned parameters  $\theta$ :

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)). \quad (4)$$

While the true posterior  $q(x_{t-1}|x_t)$  is intractable, its form conditioned on the original data  $x_0$  is tractable and Gaussian:

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t I), \quad (5)$$

with mean and variance given by:

$$\tilde{\mu}_t(x_t, x_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}x_0 + \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t, \quad (6)$$

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t. \quad (7)$$

The goal of the neural network  $\mu_\theta(x_t, t)$  is to approximate  $\tilde{\mu}_t(x_t, x_0)$ . By substituting  $x_0$  from Eq. 2 into the expression for  $\tilde{\mu}_t$ , it can be shown that predicting the mean  $\tilde{\mu}_t$  is mathematically equivalent to predicting the noise  $\epsilon$  that was added to  $x_0$  to produce  $x_t$ . This insight leads to a simplified and highly effective training objective Ho et al. (2020). Instead of optimizing the full variational lower bound (VLB), we can directly minimize the mean squared error between the true and predicted noise:

$$L_{\text{simple}}(\theta) = \mathbb{E}_{t, x_0, \epsilon} [\|\epsilon - \epsilon_\theta(x_t, t)\|^2], \quad (8)$$

where  $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$ , and  $\epsilon_\theta$  is the output of the neural network. This objective is not only computationally efficient but has also been shown to produce superior results.

## 1.4 ARCHITECTURAL FOUNDATIONS

### 1.4.1 U-NET BACKBONE

The prevalent architectural choice for many diffusion models, especially in the image domain, is the U-Net Rombach et al. (2022); Prasad et al. (2023); Xing et al. (2023). Its encoder-decoder structure with skip connections is highly effective for tasks requiring pixel-level prediction. The network takes two primary inputs: the noisy data  $x_t$  and the current timestep  $t$ . The timestep  $t$  is crucial context, as it informs the network about the level of noise it needs to remove. It is typically encoded using sinusoidal embeddings, inspired by Transformer positional encodings, and then incorporated into the network’s intermediate layers. A key enhancement in modern U-Nets is the integration of **self-attention** blocks, usually at lower-resolution feature maps, which enable the model to capture long-range spatial dependencies. For **conditional generation** (e.g., text-to-image), conditioning information  $c$  (like text embeddings from a CLIP model) is integrated into the U-Net, often via cross-attention mechanisms, allowing external signals to guide the denoising process.

### 1.4.2 DIFFUSION TRANSFORMERS (DiT)

A more recent innovation is the Diffusion Transformer (DiT) (Peebles & Xie (2023)), which replaces the convolutional U-Net with a Transformer backbone. This approach treats the diffusion process as a sequence modeling problem. A noisy image  $x_t$  is first broken down into a sequence of non-overlapping patches, which are then linearly projected to form input tokens. These tokens, along with timestep and conditional embeddings, are processed by a series of Transformer blocks. The final output tokens are decoded to predict the output noise for each patch. The primary advantage of the DiT architecture is its demonstrated **scalability**; its performance improves predictably and significantly with increased model size and computational resources, setting new benchmarks in high-fidelity image generation.

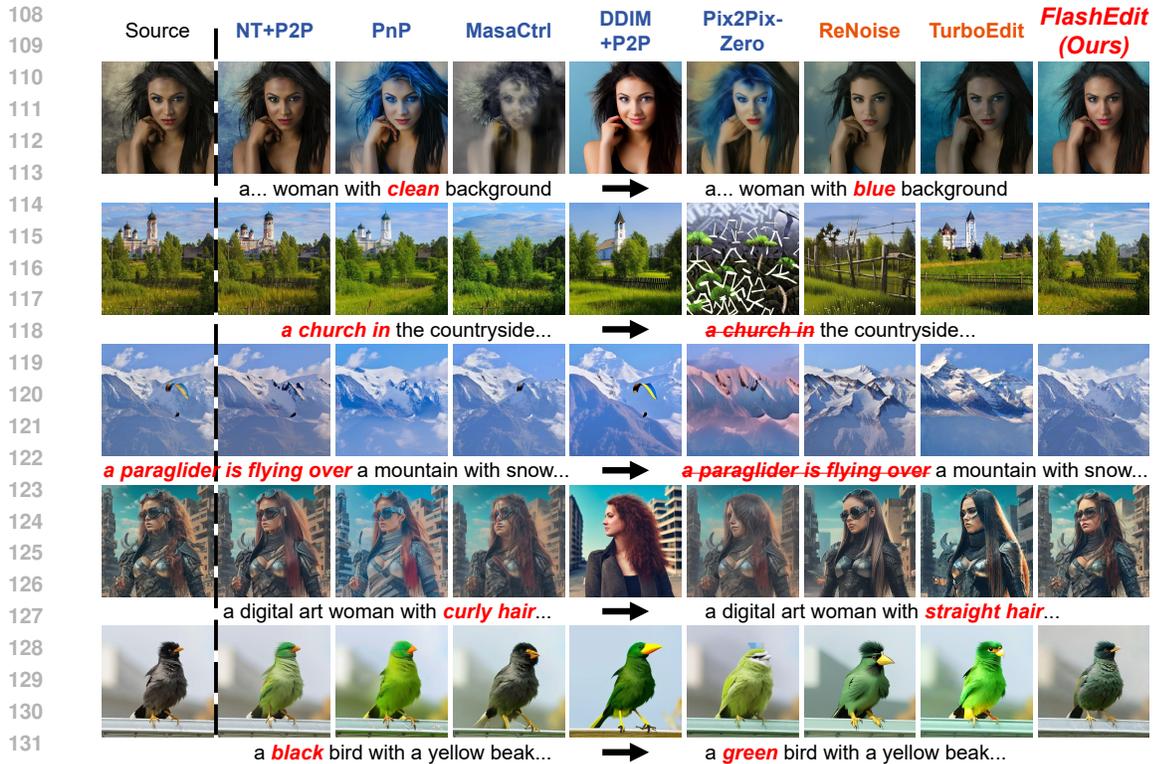


Figure 1: More comparative visual results.

## 2 EVALUATION METRICS

To quantitatively evaluate our method, we utilize a comprehensive suite of metrics designed to assess different facets of image quality, including pixel-level accuracy, perceptual similarity, and semantic consistency with the text prompt.

### 2.1 PIXEL-LEVEL FIDELITY METRICS

These metrics measure the exactness of the reconstruction at the pixel level by directly comparing the pixel values between the original and edited images.

- **Mean Squared Error (MSE)** computes the average squared difference between the pixel values of two images. It provides a strict, quantitative measure of reconstruction error. A **lower** MSE value signifies a smaller pixel-wise deviation and thus a higher-fidelity reconstruction.
- **Peak Signal-to-Noise Ratio (PSNR)** (Huynh-Thu & Ghanbari (2008)) is derived from the MSE and represents the ratio between the maximum possible pixel value (signal) and the magnitude of the reconstruction error (noise). It is expressed on a logarithmic scale (dB). A **higher** PSNR value indicates a higher quality reconstruction with less error.

### 2.2 PERCEPTUAL SIMILARITY METRICS

These metrics aim to approximate human perception of image similarity more closely than pixel-level metrics.

- **Structural Similarity Index Measure (SSIM)** (Wang et al. (2004)) evaluates image similarity by comparing three key features: luminance, contrast, and structure. It is designed to be more robust to minor pixel variations and better align with how humans perceive visual

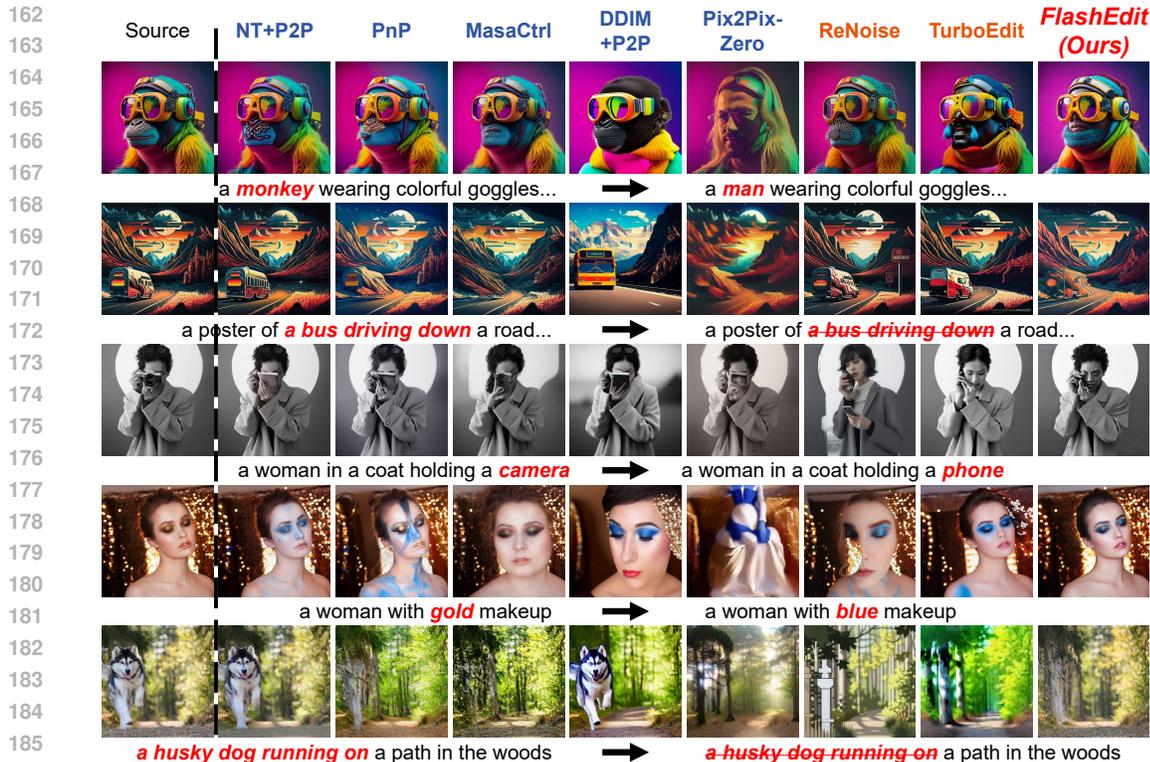


Figure 2: More comparative visual results.

quality. SSIM values range from -1 to 1, where a **higher** value closer to 1 indicates greater perceptual similarity.

- **Learned Perceptual Image Patch Similarity (LPIPS)** (Zhang et al. (2018)) measures the perceptual distance between two images using features extracted from a pre-trained deep neural network (e.g., AlexNet or VGG). By comparing images in a feature space that is more aligned with human vision, LPIPS has shown a strong correlation with human judgment of image similarity. A **lower** LPIPS score indicates that two images are more perceptually similar.

### 2.3 SEMANTIC ALIGNMENT METRIC

This metric assesses how well the generated or edited image aligns with the meaning of the input text prompt.

- **CLIP Score** (Radford et al. (2021)) measures the semantic correspondence between an image and a text description using a pre-trained CLIP (Contrastive Language-Image Pre-training) model. It calculates the cosine similarity between the image and text embeddings in a shared multimodal space. A **higher** CLIP Score indicates a stronger semantic alignment between the image content and the textual instruction.

## 3 MORE VISUAL EXPERIMENTS

To provide further evidence of our method’s efficacy and robustness, we present a comprehensive collection of additional qualitative results in Figure 1, Figure 2 and Figure 3. This section showcases the performance of **FlashEdit** across a diverse array of text-guided image editing tasks, including attribute modification, object deletion, and complex scene manipulation. We compare our results against a range of state-of-the-art multi-step and few-step baselines. The following visual examples demonstrate FlashEdit’s ability to achieve a superior balance of fidelity and precision, consistently

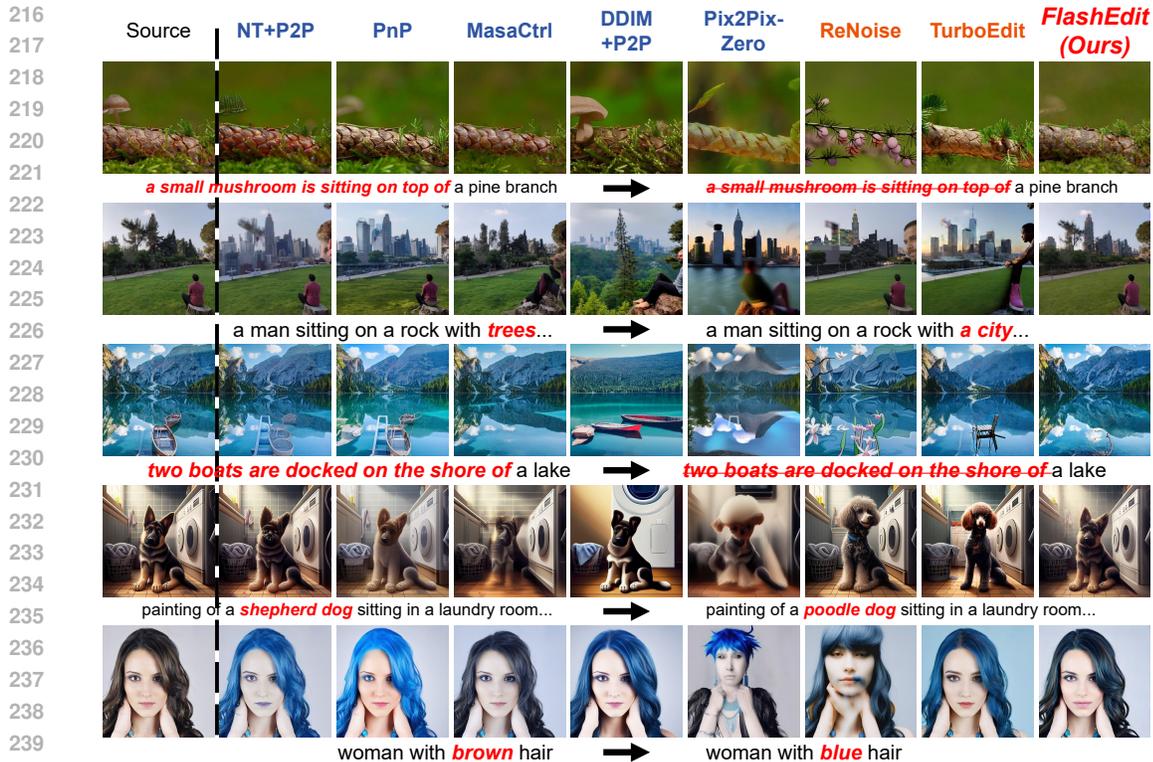


Figure 3: More comparative visual results.

243 performing localized edits with high semantic accuracy while maintaining perfect background integrity. These results visually reinforce the quantitative findings of our main paper, confirming that FlashEdit is a robust and highly effective solution for real-time generative editing.

## 247 REFERENCES

- 249 Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. Analog bits: Generating discrete data using diffusion models with self-conditioning. In *ICLR*, 2023.
- 251 Ali Hatamizadeh, Jiaming Song, Guilin Liu, Jan Kautz, and Arash Vahdat. Diffit: Diffusion vision transformers for image generation. In *ECCV*, 2024.
- 254 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- 256 Q. Huynh-Thu and M. Ghanbari. Scope of validity of psnr in image/video quality assessment. *Electronics Letters*, 44:800–801, 2008. doi: 10.1049/el:20080522. URL <https://digital-library.theiet.org/doi/abs/10.1049/el%3A20080522>.
- 260 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *NeurIPS*, 2022.
- 263 Gwanghyun Kim, Alonso Martinez, Yu-Chuan Su, Brendan Jou, José Lezama, Agrim Gupta, Lijun Yu, Lu Jiang, Aren Jansen, Jacob Walker, et al. A versatile diffusion transformer with mixture of noise levels for audiovisual generation. *arXiv preprint arXiv:2405.13762*, 2024.
- 267 Diederik P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. On density estimation with diffusion models. In *NeurIPS*, 2021.
- 269 Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.

- 270 Zhenghao Lin, Yeyun Gong, Yelong Shen, Tong Wu, Zhihao Fan, Chen Lin, Weizhu Chen, and Nan  
271 Duan. Genie: Large scale pre-training for text generation with diffusion model. *arXiv preprint*  
272 *arXiv:2212.11685*, 2022.
- 273  
274 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.  
275 In *ICCV*, 2021.
- 276 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*  
277 *the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.
- 278  
279 Vidya Prasad, Chen Zhu-Tian, Anna Vilanova, Hanspeter Pfister, Nicola Pezzotti, and Hendrik  
280 Strobelt. Unraveling the temporal dynamics of the unet in diffusion models. *arXiv preprint*  
281 *arXiv:2312.14965*, 2023.
- 282 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
283 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
284 models from natural language supervision. In *International conference on machine learning*, pp.  
285 8748–8763. PmLR, 2021.
- 286 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
287 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*  
288 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 289  
290 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *ArXiv*,  
291 *abs/2202.00512*, 2022.
- 292  
293 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*,  
294 2021.
- 295 Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error  
296 visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.  
297 doi: 10.1109/TIP.2003.819861.
- 298  
299 Zhaohu Xing, Liang Wan, Huazhu Fu, Guang Yang, and Lei Zhu. Diff-unet: A diffusion embedded  
300 network for volumetric segmentation. In *MICCAI*, 2023.
- 301  
302 Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang,  
303 Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models  
with an expert transformer. In *ICLR*, 2025.
- 304  
305 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable  
306 effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- 307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323