

EEdit : Rethinking the Spatial and Temporal Redundancy for Efficient Image Editing

Anonymous ICCV submission

Paper ID 4183

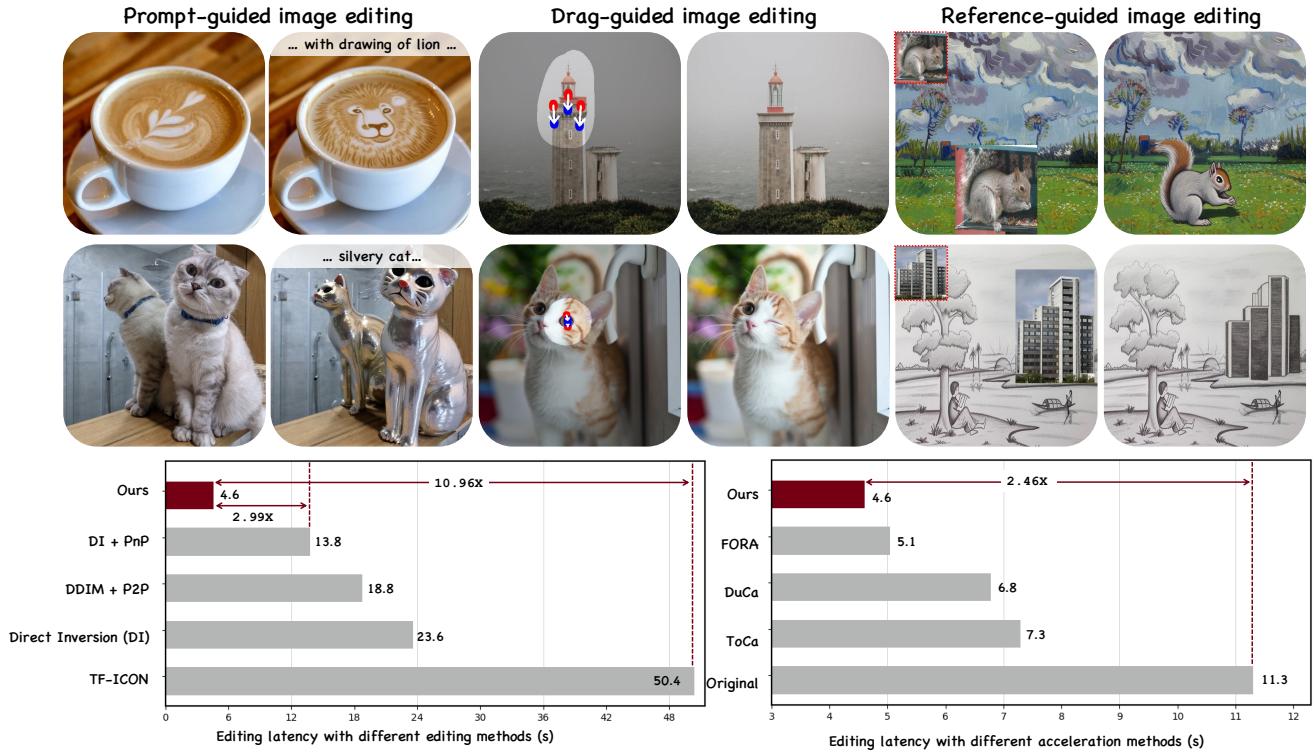


Figure 1. **Gallery of various editing tasks and efficiency comparisons.** We propose the EEdit, a novel inversion-based framework for efficient image editing. Compare with previous methods, we achieve the faster and more efficient image editing.

Abstract

001 *Inversion-based image editing is rapidly gaining momentum while suffering from significant computation overhead,*
 002 *hindering its application in real-time interactive scenarios.*
 003 *In this paper, we rethink that the redundancy in inversion-*
 004 *based image editing exists in both the spatial and tem-*
 005 *poral dimensions, such as the unnecessary computation*
 006 *in unedited regions and the redundancy in the inversion*
 007 *progress. To tackle these challenges, we propose a prac-*
 008 *tical framework, named **EEdit**, to achieve efficient image*
 009 *editing. Specifically, we introduce three techniques to solve*
 010

them one by one. **For spatial redundancy**, spatial locality caching is introduced to compute the edited region and its neighboring regions while skipping the unedited regions, and token indexing preprocessing is designed to further accelerate the caching. **For temporal redundancy**, inversion step skipping is proposed to reuse the latent for efficient editing. Our experiments demonstrate an average of **2.46**× acceleration without performance drop in a wide range of editing tasks including prompt-guided image editing, dragging and image composition. Our codes are available in the supplementary material and will be released on Github.

011
012
013
014
015
016
017
018
019
020
021

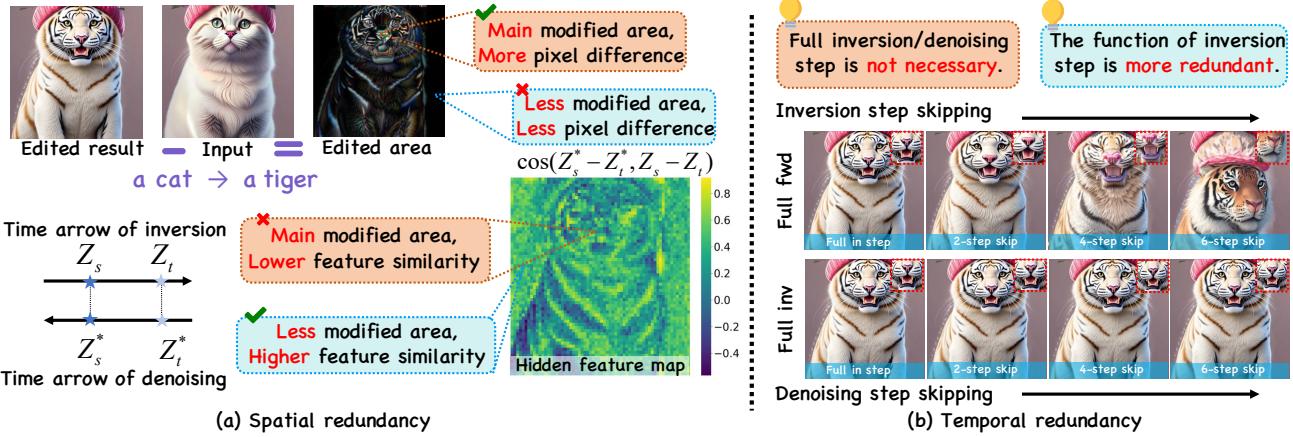


Figure 2. **The illustration of motivation.** We observe that there are significant computational redundancy during the editing process, which can be classified into spatial redundancy (a) and temporal redundancy (b). For spatial redundancy (a), due to a high background similarity between the source image and the edited results, calculating the entire image in every inversion/denoising step is unnecessary. For temporal redundancy (b), we discover that skipping some steps in inversion/denoising does not degrade the editing performance.

022

1. Introduction

023 With groundbreaking advances in diffusion models [6, 8, 024 30, 44], they have emerged as the state-of-the-art approach 025 for both image generation and editing tasks [4, 11, 13, 45]. 026 Current diffusion-based editing frameworks typically adopt 027 a two-stage pipeline: (1) an *inversion process* [27, 29, 33, 028 41] that maps the input image to its corresponding noisy latent, 029 followed by (2) a *denoising process* [16, 24, 33, 38] that 030 progressively denoises and modifies the latent code 031 to produce the edited image. While achieving impressive 032 editing quality, this two-stage paradigm suffers from 033 significant computational overhead, particularly hindering its 034 usage from real-time interactive applications on resource- 035 constrained edge devices. To address this problem, this 036 paper begins by identifying two types of computational redundancy 037 in editing and then solves them one by one.

038 **(I) Spatial Redundancy:** Although only a small region 039 of the image is expected to be edited, the current editing 040 pipeline has to compute all the pixels of the given image. 041 Retaining computation in these unedited regions incurs ad- 042 dditional computational overhead, while yielding minimal 043 benefits for our editing objectives. In Figure. 2, we visualize 044 and compare the differences between pixels before and after 045 image editing, as well as the heat map of cosine similarity 046 of the latent space state differences at fixed time intervals. 047 Both patterns consistently exhibit spatial redundancy dif- 048 ferences in the editing task: the edited regions show greater 049 pixel differences and lower similarity in the latent space, 050 whereas the unedited regions exhibit negligible pixel dif- 051 ferences and higher similarity in the latent space. Specif- 052 ically, the unedited region exhibits significantly greater re- 053 dundancy compared to the edited region.

054 **(II) Temporal Redundancy:** Compared with generating 055 a new image, image editing poses additional computation 056 cost for inversion, which is employed to map the to-be- 057 edited image to the latent space of noise to some extent. 058 In the traditional editing pipeline, the inversion process 059 takes the same computation cost as denoising and thus dou- 060 bles the computational overhead. In Figure. 2, we sepa- 061 rately control the full inversion and denoising processes by 062 introducing interval-skipping time steps during the inver- 063 sion/denoising process. It allows us to compare how reduc- 064 ing the computational load of one process under the same 065 overall computational cost affects the editing performance. 066 Notably, reducing the denoising steps first led to the loss of 067 fine details, such as texture information around the mouth 068 and nose, and soon resulted in the degradation of overall 069 structural integrity. Surprisingly, skipping time steps in the 070 inversion process had little to no perceptible effect, indi- 071 cating a highly unbalanced distribution of temporal redun- 072 dancy between inversion and denoising. Specifically, we 073 discover that inversion exhibits significantly greater redun- 074 dancy compared to denoising.

075 Based on this observation, we propose **Spatial Locality** 076 **Caching (SLoC)**, which aims to skip the computation of 077 unedited regions by reusing their features computed in the 078 previous timesteps. Concretely, during both the denois- 079 ing and inversion process, SLoC first computes the 080 tokens corresponding to all the regions and stores them in 081 a cache. Then, in the following step, SLoC still performs 082 full computation on the edited region and its neighboring 083 regions, while skipping the computation of other regions 084 by reusing their previously cached features. Such a mixed- 085 computation manner in SLoC enables diffusion models to 086 pay more effort in important regions guided by the editing

087 prior, maintaining good quality in the edited region while
088 trading the computation in unedited region for efficiency.
089 Furthermore, to further reduce the computation from SLoC,
090 we analyze that caching initialization and update strategy
091 for the score map can be finished as an offline operation.
092 Based on this observation, we propose the token index pre-
093 processing, achieving over a **15%** improvement in infer-
094 ence speed. Importantly, this optimization is mathemati-
095 cally equivalent, ensuring that SLoC itself undergoes an ad-
096 dditional lossless acceleration.

097 Additionally, we observe that there is temporal redun-
098 dancy in both inversion/denoising processing. Motivated by
099 DDIM [38], we propose the inversion step skipping strat-
100 egic. The primary insight is that *skipping certain inversion*
101 *steps does not produce noticeable artifacts and can*
102 *markedly improve processing speed*. Surprisingly, our ex-
103 periments reveal that the number of timesteps for inversion
104 can be safely reduced to 33.3% of that for diffusion, with
105 almost no noticeable performance degradation.

106 We validated the effectiveness of proposed modules in
107 extensive experiments on various editing tasks, including
108 prompt-guided [6, 13, 42], reference-guided [24, 43, 46],
109 and drag-guided editing [22, 28, 36, 48]. While achiev-
110 ing state-of-the-art performance in background consistency,
111 EEdit also delivers up to **10.96 \times** latency acceleration com-
112 pared to other editing methods. Even when compared
113 to existing cache-based acceleration techniques, it demon-
114 strates a superior acceleration ratio. Furthermore, EEdit
115 achieves a **2.46 \times** speedup with minimal performance dif-
116 ferences, maintaining near-lossless quality across various
117 evaluation metrics when compared to the original editing
118 pipeline without cache acceleration.

119 In summary, our contributions include:

- 120 • We emphasize the spatial and temporal redundancy in
121 inversion-based editing tasks and propose EEdit, a novel
122 editing framework to modify images efficiently.
- 123 • **Spatial Locality Caching (SLoC):** To reduce the spatial
124 redundancy in editing, spatial locality caching is proposed
125 to skip most of the computation in unedited regions. Then
126 we design the **Token Index Preprocessing (TIP)** to fur-
127 ther optimize the speed of caching.
- 128 • **Inversion Step Skipping (ISS):** To reduce the tempo-
129 ral redundancy, we propose to assign more computation
130 to denoising and fewer computation to inversion, which
131 firstly demonstrates and leverages the unequal importance
132 of the two processes to accelerate editing.
- 133 • **Extensive Adaption and Experiments:** We evaluate our
134 methods across various editing tasks, including prompt-
135 guided editing, drag-guided editing, and reference-guided
136 editing, which demonstrates 10.96 \times acceleration than the
137 current state-of-the-art editing approach.

2. Related Work

2.1. Acceleration of Diffusion Models

140 Low-latency, high-quality generation methods are an im-
141 portant research direction. Currently, there are two main
142 types of diffusion model acceleration methods: the first one
143 is to reduce the number of sampling steps [19–21, 38], and
144 the second one is to accelerate internal computation of the
145 diffusion model. Solutions to reduce computational com-
146 plexity include model distillation and compression [44],
147 token pruning [47], token merging [2, 3], and layer-wise
148 caching techniques [18, 25, 26, 35]. However, layer-wise
149 caching techniques have a large cache granularity, which
150 ignore the asymmetry of importance at the token level.
151 Toca [50] and Duca [51] adopt a token-wise cache, focus-
152 ing on and assigning importance to tokens with score maps.
153 During recomputation, a certain proportion of tokens are se-
154 lected for refreshing, achieving a lossless acceleration.

155 Unfortunately, existing caching schemes are not specif-
156 ically designed to accommodate the characteristics of edit-
157 ing tasks and suffer from the following issues:

158 First, computing intrinsic feature correlations between
159 tokens introduces additional internal overhead in caching
160 operations. Second, some caching strategies and editing
161 methods require accessing attention maps [50] or storing
162 KV matrices [24, 49], leading to incompatibility with exist-
163 ing Transformer acceleration techniques, such as FlashAt-
164 tention [7], thereby increasing latency. Moreover, editing
165 tasks inherently exhibit prior information regarding the spa-
166 tial distribution of regions requiring greater attention, yet
167 current caching techniques fail to leverage this information.

168 Our approach significantly improves acceleration over
169 traditional caching methods by effectively reducing both
170 temporal and spatial redundancy. Additionally, we intro-
171 duce token index preprocessing to further compress the ex-
172 tra overhead incurred by caching. As a result, our method
173 simultaneously optimizes both editing performance and ac-
174 celeration efficiency among various editing tasks.

2.2. Image Editing

175 Image editing, as an important application in the generative
176 field, has received widespread research and exploration in
177 the academic community. This includes controllable text-
178 to-image generation, image inpainting, and image-to-image
179 generation, among other approaches [6, 9, 10, 14, 31, 34].
180 Common editing approaches follow a noise addition and de-
181 noising framework, where the original image is perturbed
182 by a certain level of noise in the latent space to leverage the
183 model’s editing capabilities. The final edited image is then
184 obtained through a denoising process that restores clarity.
185 Training-free inversion editing techniques, when applied to
186 editing tasks such as prompt-guided editing, image compo-
187 sition, and image dragging, involve operations on the at-

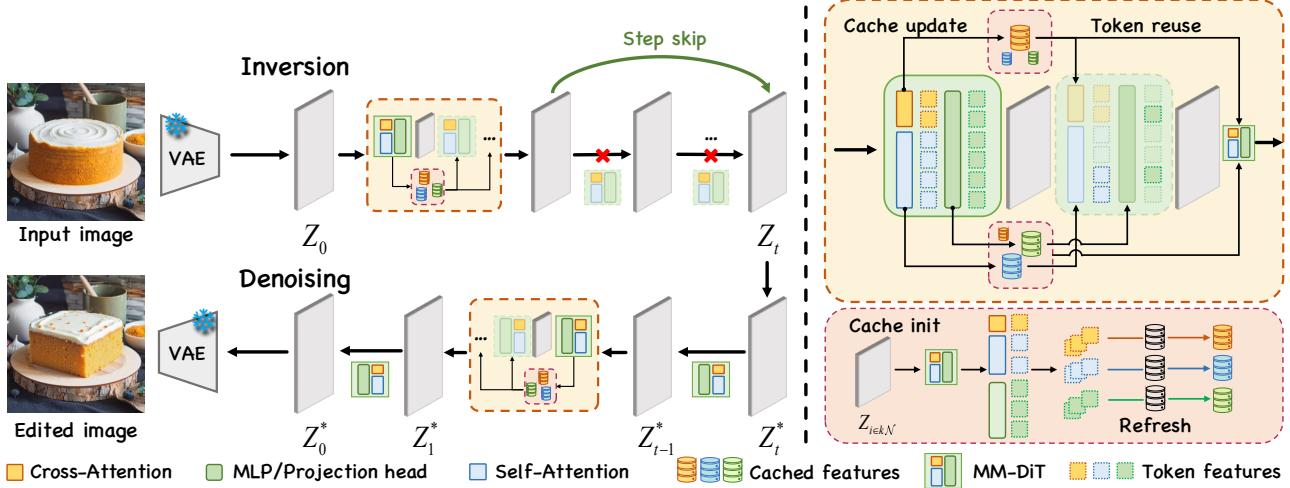


Figure 3. **The overview of our approach.** The proposed framework for image editing based on MM-DiT diffusion models employs an efficient denoising and training-free approach. The pipeline takes the original image and an editing prompt as inputs. Specifically, the cache is refreshed entirely in fixed time-step interval, while partial computation for updating cache is maintained for the intermediate timesteps.

189 tention map, including modification, enhancement, and re-
 190 placement. These methods have been applied in P2P [11]
 191 and more subsequent works [24, 49]. Since inversion-based
 192 approaches require the inversion technique to reconstruct
 193 the denoised image and significantly impact editing quality,
 194 researchers have explored both training-based and training-
 195 free inversion techniques. InfEdit adopts a virtual inver-
 196 sion strategy without explicit inversion during sampling, en-
 197 abling accurate and consistent editing, and falls under the
 198 category of inversion-free image editing methods. Unfor-
 199 tunately, these methods based on or utilizing attention maps
 200 require storing the attention map and manipulating the cor-
 201 responding KV matrices [11, 49]. This results in incompat-
 202 ibility with common attention acceleration techniques, in-
 203 creasing inference latency. In contrast, our method en-
 204 ables reduces the redundancy in inversion and denoising, enhanc-
 205 ing the efficiency and speed of editing.

206 3. Preliminaries

207 3.1. Rectified Flow

208 The Rectified Flow [17] method models the transformation
 209 from a Gaussian noise distribution π_0 to the real data dis-
 210 tribution π_1 as a continuous change along a straight path by
 211 learning a forward simulation system.

212 In the forward process, the state of the system can be
 213 viewed as a linear interpolation between the initial state and
 214 the Gaussian noise, which is simplified and easier to under-
 215 stand compared to traditional methods [12].

$$216 \mathbf{X}_t = (1-t)\mathbf{X}_1 + t\mathbf{X}_0, \quad \mathbf{X}_1 \sim \pi_1, \mathbf{X}_0 \sim \pi_0 \quad (1)$$

217 By differentiating the above expression with respect to
 218 t , the ODE form of the Rectified Flow can be obtained:

219 $\frac{d\mathbf{X}_t}{dt} = \mathbf{X}_1 - \mathbf{X}_0$. Furthermore, let us define $v_\theta = \frac{d\mathbf{X}_t}{dt}$, the
 220 training objective can then be transformed into minimizing
 221 the integral of this expectation over the time steps:

$$222 \min_\theta \int_0^1 \mathbb{E} \left[\|\mathbf{X}_1 - \mathbf{X}_0 - v_\theta(\mathbf{X}_t, t)\|^2 \right] dt \quad (2)$$

223 Here, θ represents the parameters of the neural network to
 224 be optimized. Due to the complexity introduced by the in-
 225 tegral symbol, the optimization of these parameters is typ-
 226 ically performed using the equivalent form of Conditional
 227 Flow Matching (CFM):

$$228 \mathcal{L}_{CFM} = \mathbb{E}_{t, p_t(z|\epsilon), p(\epsilon)} \|v_\theta(z, t) - u_t(z|\epsilon)\|_2^2, \quad (3)$$

229 where the conditional vector fields $u_t(z|\epsilon)$ provides an
 230 equivalent yet tractable objective. p_t is the probability path
 231 between p_0 and p_1 , and $p_1 \sim \pi_0$.

232 4. Approaches

233 Our approaches aim to reduce the spatial and temporal re-
 234 dundancy to improve the efficiency of image editing. The
 235 core idea is to leverage the mask of edited area to guide
 236 cache refreshing and reuse. The pipeline of the EEdit is
 237 shown in Figure 3. In this section, we will describe the
 238 design of spatial locality caching (see Section 4.1) and to-
 239 ken index preprocessing (see Section 4.2) for spatial redun-
 240 dancy. For temporal redundancy, we introduce the inversion
 241 step skipping in the Section 4.3.

242 4.1. Spatial Locality Caching

243 **Score Bonus for Editing Region.** In mask-guided image
 244 editing tasks, the tokens corresponding to image patches

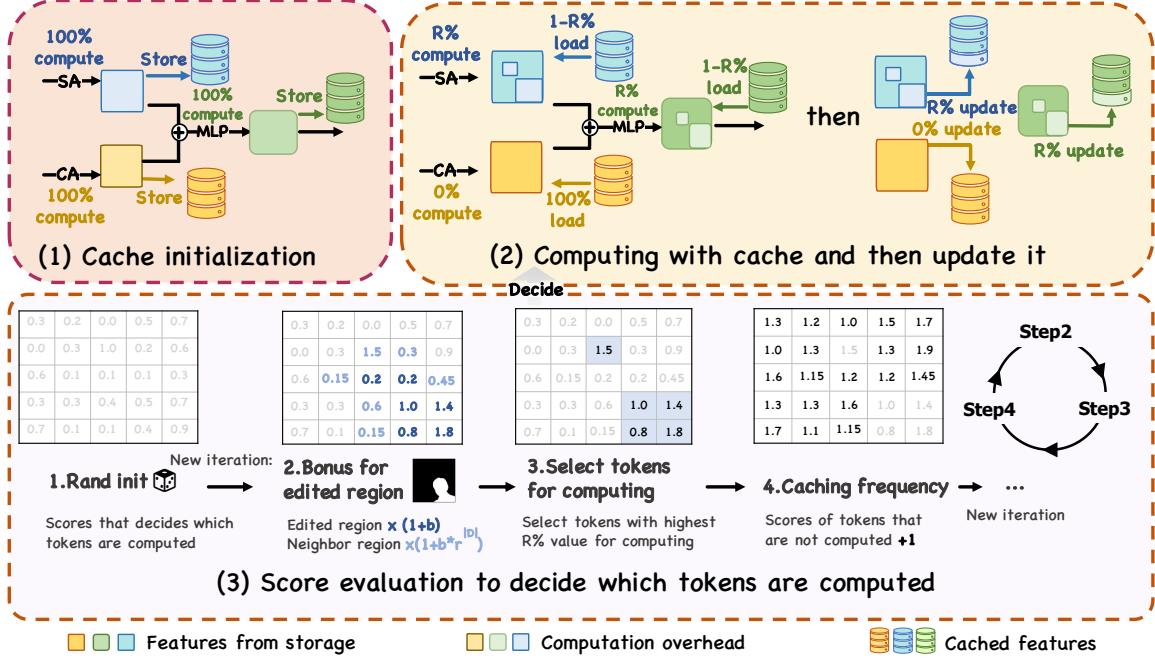


Figure 4. The pipeline of spatial locality caching. (1) The initialization and refresh process of cache storage using the computed results from SA (Self-Attention), CA (Cross-Attention), and MLP. (2) The token-wise partial computation logic and the cache update mechanism. (3) The initialization and update logic for scoring, which is responsible for selecting indices for partial computation.

outside the edited region are actually derived from Z_T^{inv} , which is obtained from inversion process. This part is directly replaced in the latent space.

As a result, the importance of DiT’s computation on these tokens is significantly reduced. Instead, we aim to focus more on the computation of tokens within the edited region. Simultaneously, to account for the modeling of correlations between important tokens and their neighboring tokens, the neighboring tokens are rewarded with gradually decreasing importance. We designed a score bonus map S_E to control the reward intensity, defined as:

$$S_E(x) = \begin{cases} 1 + b \cdot r^k, & x \in \mathcal{N}_k(M_s), k \in 0, 1, \dots, K \\ 1, & x \notin \bigcup_{k=0}^K \mathcal{N}_k(M_s) \end{cases}$$

Here x represents an arbitrary token in DiT. M_s is the latent mask for point set of edit region. Hyperparameters such as $b > 1$ represent the bonus factor, while $0 < r < 1$ denotes the decay ratio. K is the maximum number of layers for the neighborhood. The neighborhood of M_s with an L1 norm equal to k is denoted as $\mathcal{N}_k(M_s)$ and is defined as:

$$\mathcal{N}_k(M_s) = \{x : \exists e \in M_s, \|x - e\|_1 = k\}.$$

Update Strategy with Cache Frequency Control. When the proportion of tokens corresponding to the edited region is lower than the full computation of the network, meaning it is smaller than the cache refresh ratio, not all tokens

in the edited region undergo the same level of full computation. Instead, we prioritize recomputing and refreshing in the cache those tokens that have been updated less frequently and reused more times.

SLoC track the frequency of times each token is reused from the cache and increment its score accordingly on a token-wise map M_{freq} . The more frequently a token is reused, the more likely its features will be recomputed and refreshed. Once certain tokens undergo recomputation and refresh, their corresponding usage frequency counters are immediately reset to zero. From another perspective, the design of cache frequency control serves two key purposes. First, it encourages the recomputation of frequently reused tokens to reduce accumulated errors. Second, it suppresses the redundant recomputation of tokens that are repeatedly updated, thereby reducing computational overhead.

4.2. Token Index Preprocessing

Although SLoC reduces spatial redundancy by using caching, as a double-edged sword, such an acceleration comes with certain limitations: in a cache cycle shown in Figure. 4, the cache overhead is sequentially arranged as follows: (1) Randomly initialize. (2) Bonus for edited region. (3) Perform sorting and selection. (4) Update. (5) Refresh. These steps constitute the caching cycle and operate as illustrated in Figure. 4.

However, we observe that internal score updates and token selection take additional computational costs in editing

Algorithm 1 SLoC Editing with ISS&TIP

Require: Input image \mathbf{I}_s , Mask for editing region \mathbf{M}_s ,
 Prompt for editing \mathbf{P}_m , Randomly initialized map \mathcal{R} ,
 Bonus for edited region \mathbf{S}_E and Cache dict $\mathcal{C}_{l,m}[:, :]$.

Ensure: The edited result \mathbf{I}^*

- 1: // **Token Index Preprocessing**
- 2: $\mathcal{M}_{freq} \leftarrow zero[\mathcal{F}_i, 1]$
- 3: **for** $t = T, T - 1, \dots, 1$ **do**
- 4: **for** $\mathcal{F}_i \leftarrow \mathbf{SA}_l, \mathbf{CA}_l, \mathbf{MLP}_l, l \in [1, 2 \dots \mathcal{L}]$ **do**
- 5: $\mathcal{S}_l \leftarrow (\mathcal{R} \odot \mathbf{S}_E) \oplus \mathcal{M}_{freq}$
- 6: $\mathcal{I}_{i,l,t} \leftarrow \text{Sel}_{topR\%}(\mathcal{S}_l)$
- 7: Update(\mathcal{M}_{freq})
- 8: **end for**
- 9: **end for**
- 10: $\mathbf{Z}_0 \leftarrow \text{cat}(\text{VQ-Encoder}(\mathbf{I}_s), \text{Txt-Encoder}(\mathbf{P}_m))$
- 11: // **Inversion Reduction**
- 12: **for** $t = 1, 2 \dots, T$ **do**
- 13: $\mathbf{Z}_t \leftarrow \begin{cases} \mathbf{Z}_{t-1} & \text{if } t \bmod m \neq 1 \text{ and } m \neq T, m \in \mathcal{N} \\ \text{RF-inversion}(\mathbf{Z}_{t-1}, t-1, \phi) & \text{otherwise} \end{cases}$
- 14: **end for**
- 15: $\mathbf{Z}_T^* \leftarrow \mathbf{Z}_T$
- 16: // **Image Editing Steps with SLoC**
- 17: **for** $t = T, T - 1, \dots, 1$ **do**
- 18: **for** $\mathcal{F}_i \leftarrow \mathbf{SA}_l, \mathbf{CA}_l, \mathbf{MLP}_l, l \in [1, \dots, \mathcal{L}]$ **do**
- 19: $\mathbf{Z}_{l+1}^* \leftarrow \text{scatter}(\mathcal{F}_i(\mathbf{Z}_l^*, \mathcal{I}_{i,l,t}), \mathcal{C}_{t+1}[l, \mathcal{F}_i])$
- 20: Update($\mathcal{C}_{t+1}[l, \mathcal{F}_i]$)
- 21: **end for**
- 22: $\mathbf{Z}_{t-1}^* \leftarrow \mathbf{Z}_{t-1}^* \odot \mathbf{M}_s + \mathbf{Z}_t \odot (1 - \mathbf{M}_s)$
- 23: **end for**
- 24: $\mathbf{I}^* \leftarrow \text{VQ-Decoder}(\mathbf{Z}_0^*)$
- 25: **return** \mathbf{I}^*

295 processing and can be further optimized. Our initialization
 296 and update strategy for the score map can be transformed
 297 from an online operation into an offline algorithm while
 298 maintaining full mathematical equivalence (See the supple-
 299 mentary material for the formal proof). The key insight here
 300 is that under the score update rule

$$301 \quad \mathcal{S} \leftarrow (\mathcal{R} \odot \mathbf{S}_E) \oplus \mathcal{M}_{freq},$$

302 we can prove the top- $R\%$ selected indices $\mathcal{I}_{topR\%}^{(t)}$ in the
 303 same time step in the offline and online process remain
 304 equivalent:

$$305 \quad \mathcal{I}_{topR\%}^{(t)}(\text{offline}) = \mathcal{I}_{topR\%}^{(t)}(\text{online}) \quad \forall t \in [1 \dots T].$$

306 Here \mathcal{S} denotes the score map, \mathcal{R} represents the randomly
 307 initialized score values, \mathbf{S}_E corresponds to the region score
 308 bonus, and \mathcal{M}_{freq} denotes the matrix for cache frequency
 309 control, consistent with Algorithm 1.

310 Therefore, we can decouple the cache score update, sort-
 311 ing, and index selection logic from the model's computation

process. By precomputing and storing the required token indices during preprocessing, the overhead from cache operations in the inference phase is reduced to a single read/write cost. Consequently, the former 4 steps of the cache cycle can be omitted during inference. SLoC directly updates only the necessary tokens during inference, without executing any redundant score computation or update strategies.

4.3. Inversion Step Skipping

After addressing the spatial redundancy, we focus on the temporal redundancy in both inversion/denoising processing. As demonstrated in Fig. 2(b), previous methods [11, 33] calculate the noise at each timestep. However, the emergence of DDIM [38] motivates us to consider *whether we can skip certain steps in flow matching-based inversion*. During the editing process, we discover that skipping some inversion steps does not degrade the quality of the editing results(see Fig. 2(b)), which demonstrates the redundancy of the inversion process.

To eliminate this redundancy, we propose the inversion step skipping strategy in the editing processing. The key insight of inversion step skipping is that *skipping some inversion steps does not result in noticeable artifacts, but it can significantly improve speed*. Formally, during the inversion, the skipping step is set m . We add the noise from the \mathbf{Z}_0 . In every m steps, we perform an rf-inversion [33], while for the other timesteps, we directly reuse the noise from the previous timestep. This process can be formally expressed as follows:

$$340 \quad \mathbf{Z}_t \leftarrow \begin{cases} \mathbf{Z}_{t-1} & \text{if } t \bmod m \neq 1 \text{ and } m \neq T, m \in \mathcal{N} \\ \text{RF-inversion}(\mathbf{Z}_{t-1}, t-1, \phi) & \text{otherwise} \end{cases}$$

341 Additionally, we provide the pseudocode implementa-
 342 tion in Algorithm 1. Note that, to ensure noise quality, the
 343 final step of the inversion process is always fully computed,
 344 whereas the intermediate steps employ a mixed strategy of
 345 full computation and cache-accelerated computation. In our
 346 experiment, skipping step m is set 3 to balance the quality
 347 and speed. We also provide the ablation study about the
 348 inversion skipping step in Section 5.3.

5. Experiments

5.1. Implementation Details

351 In our experiment, the base model adopted the currently
 352 popular flow matching model, FLUX-Dev [15], which con-
 353 sists of 12B parameters. Models for other qualitative results
 354 are implemented using SD series [1, 5, 32, 37] with original
 355 codebase. Our inference and editing pipeline is built upon
 356 the codebase from Hugging Face. The inversion and deno-
 357 sing step is set to 28. RF-inversion relative hyper-parameters
 358 follow original implementations [33]. We employ the text-
 359 guidance ratio of 7.0. We selected the PIE-Bench [39]
 360 Benchmark as the dataset for prompt-guided editing, the

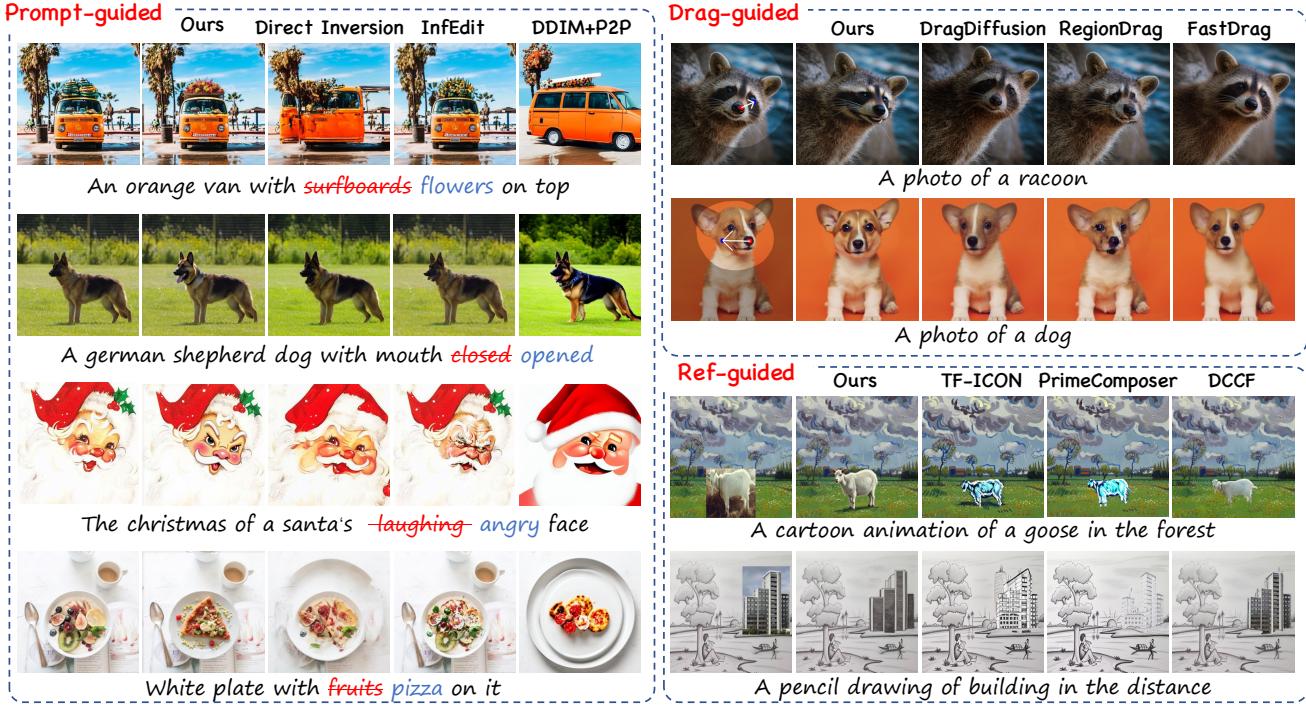


Figure 5. **Qualitative comparisons between baselines and our approach.** We compare our approach with various editing tasks, including prompt-guided image editing, drag-guided image editing, and reference-guided image editing

361 TF-ICON Test Benchmark [23] as the dataset for reference-
 362 guided editing, and DragBench-SR and DragBench-DR as
 363 the datasets [40] for drag-guided editing. All experiments
 364 were conducted on an NVIDIA H20. More details and eval-
 365 uation metrics can be found in supplementary material.

366 5.2. Compare with Baseline

367 **Qualitative Comparison.** We perform extensive qual-
 368 iitative comparisons between our method and various edit-
 369 ing approaches. Specifically, we evaluate three kinds
 370 of popular editing tasks, including prompt-guided, drag-
 371 guided, and reference-guided image editing. (1) **Prompt-
 372 guided image editing:** we compare with Direct Inversion
 373 [13], InfEdit [45], and P2P [11]. (2) **Drag-guided im-
 374 age editing,** we select three methods to do comparisons,
 375 including DragDiffusion [36], Region Drag[22], and Fast-
 376 Drag [48]. Pipeline are implemented from public codebase
 377 from github. (3) **Reference-guided editing,** we compare
 378 our approach with TF-ICON [24], PrimeComposer [43],
 379 and DCCF [46]. SD-v2-1 checkpoint is used for Image
 380 Composition. As shown in Figure. 5, for prompt-guided
 381 task, Direct inversion and InfEdit fail to edit the image suc-
 382 cessfully. DDIM+P2P has a challenge to maintain the back-
 383 ground consistency. By contrast, our approach exhibits su-
 384 perior consistency, enhanced detail preservation, and im-
 385 proved aesthetic quality in terms of style.

386 **Quantitative Comparison.** We compared various editing

Table 1. **Comparison of quality** across different editing methods.

Editing methods	Inversion type	PSNR \uparrow	LPIPS $\downarrow \times 10^{-2}$	SSIM \uparrow	CLIP-T \uparrow
P2P [11]	DDIM	17.87	20.88	0.72	25.13
MasaCtrl [4]	DDIM	22.19	10.54	0.80	24.02
	DI	22.69	8.73	0.82	24.39
PnP [13]	DDIM	25.23	11.27	0.80	25.42
	DI	22.46	10.55	0.80	25.48
InfEdit [45] DiT4Edit [9]	Inversion-free	28.11	5.61	0.85	25.86
	DPM-Solver++	22.85	-	-	25.39
SLoC	RF-inversion	31.97	1.96	0.94	25.37
	ISS	31.97	1.95	0.94	25.38

387 methods in terms of background consistency in non-edited
 388 regions and adherence to prompts in the edited regions.
 389 As shown in Table. 1, since our approach allocates fewer
 390 computational resources to non-edited regions and employs
 391 masking in the latent space, we achieve significantly better
 392 performance in background consistency compared to other
 393 state-of-the-art editing methods. Furthermore, in terms
 394 of prompt adherence, the Inversion step skipping strategy
 395 achieves a CLIP score comparable to that of RF-Inversion
 396 with full computation. It is also competitive with others.

397 5.3. Ablation Study

398 Extensive quantitative and qualitative ablation experiments
 399 were conducted to analyze the impact of our methods on
 400 editing quality. We measure the similarity to the input im-

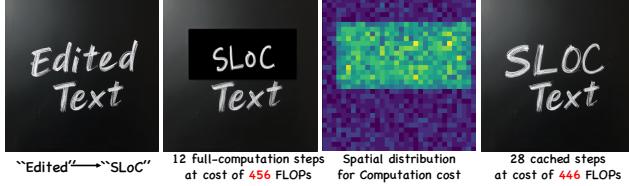


Figure 6. **A qualitative example of ablation over SLoC**, demonstrating that regional computation with SLoC leads to higher quality and lower computational overhead.

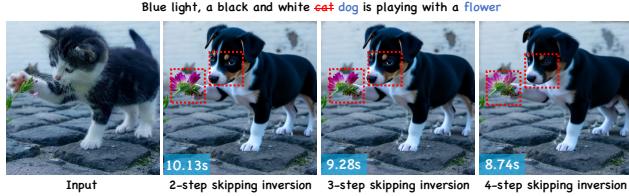


Figure 7. **The qualitative ablation study about inversion step skipping.** We visualize results in various inversion step skipping.

age in non-edited regions as *background preservation* (*BG preservation*). Additionally, we compare the editing results incorporating various modules with the editing results obtained from full computation without caching as *foreground fidelity* (*FG fidelity*).

Ablation study of spatial locality caching. Figure. 6 illustrates the impact of SLoC on editing performance. In text editing tasks, SLoC leverages the prior of spatial locality to achieve superior editing results compared to the original full-computation editing method, despite using more steps but incurring lower computational overhead. This validates the feasibility of our regionally focused computation strategy. Furthermore, quantitative comparisons with other cache-based acceleration methods demonstrate the superiority of our approach (shown in supplementary material).

Table 2. **Ablation study on ISS.** Ablation study is conducted on different skipping settings for background preservation, foreground fidelity and inference time.

Inversion	Denoising	BG preservation		FG fidelity			Inference ↓
		LPIPS $\downarrow \times 10^{-2}$	LPIPS $\downarrow \times 10^{-3}$	LPIPS $\downarrow \times 10^{-2}$	PSNR \uparrow	FID \downarrow	
Full step	Full step	1.98	-	-	-	-	13.27
2-step skip	Full step	1.98	5.46	43.77	3.35	10.16	
3-step skip	Full step	1.98	5.29	43.99	3.23	9.31	
4-step skip	Full step	1.98	5.29	43.80	3.31	8.76	

416 **Ablation study of inversion step skipping.** It can be observed in Table. 2 and Figure. 7 that increasing the skip 417 interval m to reduce inversion computation significantly 418 decreases inference time while maintaining background 419 preservation without quality degradation. Additionally, the 420 foreground fidelity metrics exhibit negligible differences 421 compared to full-step computation. This validates the 422 effectiveness and fidelity of our ISS strategy.

423 **Ablation study of TIP over ISS and various tasks.** As 424

Table 3. **Ablation study on different configurations of TIP and ISS** in prompt-guided, drag-guided, and ref-guided editing.

TASK	Method	FG fidelity			Inference (s) ↓	CLIP-E ↑		
		TIP	ISS	FID↓	PSNR↑	LPIPS $\downarrow \times 10^{-2}$		
Prompt	✗	✗	✗	39.50	31.75	5.75	5.96	21.34
	✗	✓	✗	39.33	31.76	5.74	5.06	21.34
	✓	✗	✗	39.42	31.75	5.75	5.14	21.34
	✓	✓	✗	39.21	31.76	5.74	4.60	21.34
Dragging	✗	✗	✗	20.61	33.47	2.28	7.12	22.20
	✓	✓	✗	22.07	33.68	2.19	5.66	22.21
Composition	✗	✗	✗	12.33	39.78	0.54	7.25	23.35
	✓	✓	✗	12.35	39.80	0.54	5.66	23.35

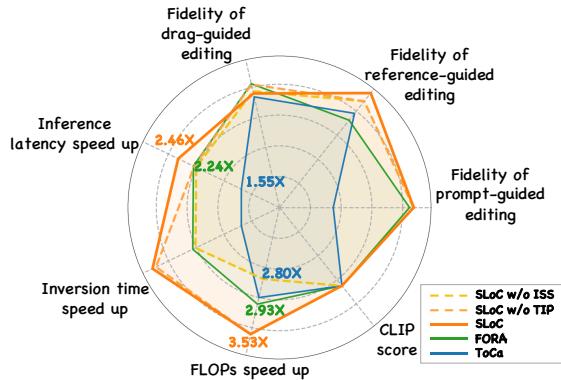


Figure 8. **A Comparison over different configurations and other cache methods.** Fidelity of various tasks, speed up ratio and clip score are shown in this radar chart.

425 shown in Table 3 and Figure 8, TIP remains compatible with 426 ISS while maintaining comparable FG fidelity and CLIP 427 scores across various editing tasks. Furthermore, it achieves 428 an additional reduction in inference latency by an average 429 of over 20% compared to the SLoC baseline. Under the 430 same configuration, the impact of incorporating TIP on editing 431 quality remains within the range of random fluctuations. 432 Furthermore, our quantitative experiments validate that this 433 approach achieves lossless acceleration while maintaining 434 the same generation quality which can be seen in supple- 435

6. Conclusions

436 Inversion-based image editing usually suffers from ex- 437 pensive computation costs caused by the spatial and 438 temporal redundancy within diffusion models. To solve 439 this problem, we design an efficient editing framework 440 using the spatial locality caching with token index pre- 441 processing and inversion step skipping. To the best of 442 our knowledge, we are the first to adapt cache-based 443 acceleration for diffusion inference across various edit- 444 ing tasks. Our work provides valuable insights and 445 exploration for future approaches toward efficient, 446 and potentially real-time, image and even video editing.

References

- [1] Stability AI. Stable diffusion 2.1 base. <https://huggingface.co/stabilityai/stable-diffusion-2-1-base>, 2022. Accessed: March 5, 2025. 6
- [2] Daniel Bolya and Judy Hoffman. Token merging for fast stable diffusion. *CVPR Workshop on Efficient Deep Learning for Computer Vision*, 2023. 3
- [3] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster, 2023. 3
- [4] Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. Masactr: Tuning-free mutual self-attention control for consistent image synthesis and editing, 2023. 2, 7
- [5] CompVis. Stable diffusion v1.4. <https://huggingface.co/CompVis/stable-diffusion-v1-4>, 2022. Accessed: March 5, 2025. 6
- [6] Y. Dalva, K. Venkatesh, and P. Yanardag. Fluxspace: Disentangled semantic editing in rectified flow transformers. 2024. 2, 3
- [7] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness, 2022. 3
- [8] 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. 2
- [9] K. Feng, Y. Ma, B. Wang, et al. Dit4edit: Diffusion transformer for image editing. 2024. 3, 7
- [10] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion, 2022. 3
- [11] A. Hertz, R. Mokady, J. Tenenbaum, K. Aberman, Y. Pritch, and D. Cohen-Or. Prompt-to-prompt image editing with cross attention control. 2022. 2, 4, 6, 7
- [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint arXiv:2006.11239*, 2020. 4
- [13] Xuan Ju, Ailing Zeng, Yuxuan Bian, Shaoteng Liu, and Qiang Xu. Pnp inversion: Boosting diffusion-based editing with 3 lines of code. *International Conference on Learning Representations (ICLR)*, 2024. 2, 3, 7
- [14] Vladimir Kulikov, Matan Kleiner, Inbar Huberman-Spiegelglas, and Tomer Michaeli. Flowedit: Inversion-free text-based editing using pre-trained flow models. *arXiv preprint arXiv:2412.08629*, 2024. 3
- [15] Black Forest Labs. Flux.1-dev. <https://huggingface.co/black-forest-labs/FLUX.1-dev>, 2024. Accessed: March 5, 2025. 6
- [16] S. Lin, B. Liu, J. Li, and X. Yang. Common diffusion noise schedules and sample steps are flawed. 2024. 2
- [17] Y. Lipman, R.T.Q. Chen, H. Ben-Hamu, M. Nickel, and M. Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2023. 4
- [18] H. Liu, W. Zhang, J. Xie, et al. Faster diffusion via temporal attention decomposition. 2024. 3
- [19] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow, 2022. 3
- [20] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps, 2022.
- [21] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models, 2023. 3
- [22] Jingyi Lu, Xinghui Li, and Kai Han. Regiondrag: Fast region-based image editing with diffusion models, 2024. 3, 7
- [23] Shilin Lu, Yanzhu Liu, and Adams Wai-Kin Kong. Tf-icon: Diffusion-based training-free cross-domain image composition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2294–2305, 2023. 7
- [24] Shilin Lu, Yanzhu Liu, and Adams Wai-Kin Kong. Tf-icon: Diffusion-based training-free cross-domain image composition, 2023. 2, 3, 4, 7
- [25] X. Ma, G. Fang, and X. Wang. Deepcache: Accelerating diffusion models for free. *arXiv preprint arXiv:2312.00858*, 2023. 3
- [26] X. Ma, G. Fang, M.B. Mi, and X. Wang. Learning-to-cache: Accelerating diffusion transformer via layer caching. *arXiv preprint arXiv:2406.01733*, 2024. 3
- [27] R. Mokady, A. Hertz, K. Aberman, Y. Pritch, and D. Cohen-Or. Null-text inversion for editing real images using guided diffusion models. 2022. 2
- [28] Shen Nie, Hanzhong Allan Guo, Cheng Lu, Yuhao Zhou, Chenyu Zheng, and Chongxuan Li. The blessing of randomness: Sde beats ode in general diffusion-based image editing, 2024. 3
- [29] Z. Pan, R. Gherardi, X. Xie, and S. Huang. Effective real image editing with accelerated iterative diffusion inversion. 2023. 2
- [30] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 2
- [31] 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. 3
- [32] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022. 6
- [33] L. Rout, Y. Chen, N. Ruiz, C. Caramanis, S. Shakkottai, and W.-S. Chu. Semantic image inversion and editing using rectified stochastic differential equations. *arXiv preprint arXiv:2410.10792*, 2024. 2, 6

- 561 [34] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch,
562 Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine
563 tuning text-to-image diffusion models for subject-driven
564 generation, 2023. 3
- 565 [35] P. Selvaraju, T. Ding, T. Chen, I. Zharkov, and L. Liang.
566 Fora: Fast-forward caching in diffusion transformer accel-
567 eration. *arXiv preprint arXiv:2407.01425*, 2024. 3, 2
- 568 [36] Yujun Shi, Chuhui Xue, Jiachun Pan, Wenqing Zhang, Vin-
569 cent YF Tan, and Song Bai. Dragdiffusion: Harnessing diffu-
570 sion models for interactive point-based image editing. *arXiv
571 preprint arXiv:2306.14435*, 2023. 3, 7
- 572 [37] SimianLuo. Lcm dreamshaper v7. https://huggingface.co/SimianLuo/LCM_Dreamshaper_v7, 2024. Accessed: March 5, 2025. 6
- 573 [38] J. Song, C. Meng, and S. Ermon. Denoising diffusion im-
574 plicit models. 2022. 2, 3, 6
- 575 [39] PIE-Bench Team. Pie-bench: Prompt-driven im-
576 age editing benchmark. <https://forms.gle/hVMkTABb4uvZVjme9>, 2024. Accessed: March 5, 2025.
577 6
- 581 [40] Visual AI Team. Regiondrag: Interactive region-based im-
582 age editing. <https://github.com/Visual-AI/RegionDrag>, 2024. Accessed: March 5, 2025. 7
- 584 [41] J. Wang, J. Pu, Z. Qi, et al. Taming rectified flow for inver-
585 sion and editing. 2024. 2
- 586 [42] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao,
587 Jan Kautz, and Bryan Catanzaro. High-resolution image
588 synthesis and semantic manipulation with conditional gans,
589 2018. 3
- 590 [43] Y. Wang, W. Zhang, J. Zheng, and C. Jin. Primecomposer:
591 Faster progressively combined diffusion for image composi-
592 tion with attention steering. 2024. 3, 7
- 593 [44] E. Xie, J. Chen, J. Chen, et al. Sana: Efficient high-resolution
594 image synthesis with linear diffusion transformers. 2024. 2,
595 3
- 596 [45] S. Xu, Y. Huang, J. Pan, Z. Ma, and J. Chai. Inversion-free
597 image editing with natural language. 2023. 2, 7
- 598 [46] Ben Xue, Shenghui Ran, Quan Chen, Rongfei Jia, Binqiang
599 Zhao, and Xing Tang. Dccf: Deep comprehensible color
600 filter learning framework for high-resolution image harmo-
601 nization, 2022. 3, 7
- 602 [47] Evelyn Zhang, Jiayi Tang, Xuefei Ning, and Linfeng Zhang.
603 Training-free and hardware-friendly acceleration for diffu-
604 sion models via similarity-based token pruning. In *Proceed-
605 ings of the AAAI Conference on Artificial Intelligence*, 2025.
606 3
- 607 [48] Xuanjia Zhao, Jian Guan, Congyi Fan, Dongli Xu, Youtian
608 Lin, Haiwei Pan, and Pengming Feng. Fastdrag: Manipulate
609 anything in one step. In *The Thirty-eighth Annual Confer-
610 ence on Neural Information Processing Systems*, 2024. 3,
611 7
- 612 [49] Tianrui Zhu, Shiyi Zhang, Jiawei Shao, and Yansong Tang.
613 Kv-edit: Training-free image editing for precise background
614 preservation. *arXiv preprint arXiv:2502.17363*, 2025. 3, 4
- 615 [50] C. Zou, X. Liu, T. Liu, S. Huang, and L. Zhang. Accelerat-
616 ing diffusion transformers with token-wise feature caching.
617 *arXiv preprint arXiv:2410.05317*, 2024. 3, 2

618
619
620