

BlazeEdit: Generalist Image Editing on Mobile Devices with Image-to-Image Diffusion Models

Anonymous CVPR submission

Paper ID 2

Abstract

001 *The remarkable generation quality of modern diffusion*
002 *models often comes at the cost of massive parameter*
003 *counts, which necessitate server-side inference with sig-*
004 *nificant computational costs and potential privacy risks.*
005 *Consequently, there is growing momentum toward devel-*
006 *oping efficient on-device alternatives. While recent efforts*
007 *have optimized text-to-image models for mobile hardware,*
008 *they remain relatively bulky, typically ranging from 0.5B*
009 *to 1B parameters. We present BlazeEdit, a highly effi-*
010 *cient, generalist image-to-image diffusion model tailored*
011 *for on-device deployment. By identifying that many practi-*
012 *cal image editing tasks do not require text-based guidance,*
013 *we eliminate the text-conditioning components and develop*
014 *a multi-task architecture that consolidates object removal,*
015 *outpainting, tone correction, relighting, and sticker gener-*
016 *ation into a single, compact model of only 195M param-*
017 *eters. BlazeEdit achieves a substantial reduction in down-*
018 *load size and memory overhead while maintaining compet-*
019 *itive generation quality. It completes a full inference pass*
020 *in just 290ms on a Pixel 10, delivering a seamless, privacy-*
021 *preserving, and lightning-fast experience for generalist im-*
022 *age editing on the edge.*

023 1. Introduction

024 Diffusion and flow-based generative models [1, 9, 14, 15,
025 20, 21] have achieved remarkable visual quality across a
026 variety of digital content creation domains. However, many
027 state-of-the-art diffusion models [4] employ heavy multi-
028 modal transformers whose parameter counts can reach up
029 to 20B [25], requiring server-side inference on high-end
030 GPUs/TPUs. This not only incurs significant computational
031 costs, but also raises potential concerns regarding the pri-
032 vacy of uploaded personal photos.

033 Consequently, there is growing interest in developing on-
034 device diffusion models that can run locally and efficiently.
035 Recent efforts, such as SnapFusion [13], SnapGen [3], and

MobileDiffusion [32], have made impressive strides in re- 036
ducing the denoiser size and latency for text-to-image gen- 037
eration. Nevertheless, these models remain relatively heavy, 038
as the text encoder increases the total model size. For in- 039
stance, SnapGen leverages multiple text encoders including 040
CLIP [17] and Gemma-2-2B [22], adding ~ 2 B parameters 041
to its 0.38B denoiser. The large model size often presents a 042
download barrier for mobile applications, where bandwidth 043
is frequently limited. 044

In this paper, we challenge the prevailing assumption 045
that a general-purpose mobile editing tool must be built 046
upon a text-to-image foundation. We observe that a sig- 047
nificant subset of common image editing actions, such as 048
removing background objects, changing the aspect ratio of 049
an image (e.g., portrait to square or landscape), and creating 050
stylized stickers from photos, can be sufficiently guided by 051
the input image plus a user-provided mask, bypassing the 052
need for text-based conditioning. 053

Building on this insight, we develop a pure image-to- 054
image framework and introduce BlazeEdit, a highly effi- 055
cient, generalist image editor tailored for mobile devices. 056
Our contributions are summarized as follows: 057

- We develop a pretraining pipeline specifically for image- 058
to-image diffusion models. By leveraging masked recon- 059
struction as the pretraining objective, we endow the base 060
model with foundational inpainting and outpainting capa- 061
bilities, facilitating data-efficient downstream finetuning. 062
- We repurpose the mask value as a universal task indica- 063
tor. Combined with a jointly trained image-and-mask en- 064
coder, this allows for simultaneous multi-task finetuning 065
that enables knowledge transfer across tasks. 066
- We achieve a single, compact model of only 195M pa- 067
rameters that consolidates five distinct editing tasks— 068
object removal, outpainting, tone correction, relighting, 069
and sticker generation. With an inference latency of just 070
290ms on a Pixel 10, our model delivers a highly interac- 071
tive user experience. 072

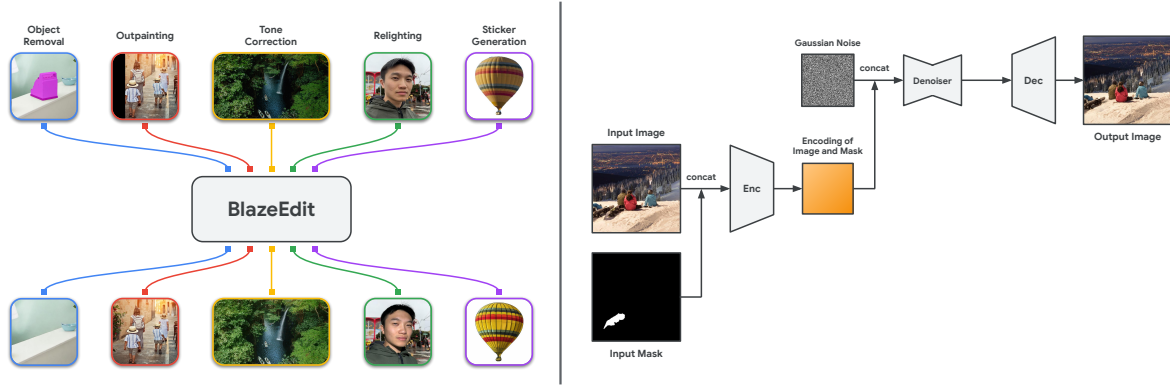


Figure 1. (Left) BlazeEdit is a generalist image-to-image editing framework designed for mobile devices. We achieve a single, compact model that consolidates five distinct editing tasks. (Right) BlazeEdit follows the latent diffusion paradigm, and employs a trainable image-and-mask encoder to provide the conditioning signal to the denoiser. The input resolution is 512×512 , and latent resolution is 64×64 .

073 2. Related Work

074 There have been significant recent advancements in devel-
 075 oping generalist text-to-image diffusion models that unify
 076 diverse, multi-modal editing tasks within single architec-
 077 tures [5, 12, 16, 25–28, 31]. However, their reliance
 078 on heavy multi-modal transformers and massive parameter
 079 counts renders them impractical for on-device execution.
 080 A growing body of research has focused on compressing
 081 text-to-image diffusion models for mobile devices. Pioneering
 082 works such as SnapFusion [13], SnapGen [3], and Mo-
 083 bileDiffusion [32] have significantly reduced the model pa-
 084 rameters to the 0.5B \sim 1B range through extensive architec-
 085 tural investigation and advanced distillation. Our work fur-
 086 ther shrinks the model size by shifting to an image-to-image
 087 framework that still supports a wide variety of practical
 088 editing tasks. Compared to the pixel-space image-to-image
 089 framework presented in Palette [19], our model works in la-
 090 tent space and introduces a task-agnostic pretraining phase,
 091 enabling fast inference and data-efficient finetuning.

092 3. Method

093 3.1. Model Architecture

094 Figure 1 presents an overview of the BlazeEdit framework.
 095 We detail our design below.

096 **Latent Diffusion with Jointly Trained Encoder.** Follow-
 097 ing established practice, we adopt latent diffusion [18]
 098 to reduce the computational overhead of the iterative denoising
 099 process. Given a frozen encoder \mathcal{E} and decoder \mathcal{D} , we
 100 seek to train a denoiser ϵ_θ that models the conditional distri-
 101 bution $p(\mathcal{E}(\mathbf{y}) \mid \mathbf{x}, \mathbf{m})$, where \mathbf{x} , \mathbf{y} , \mathbf{m} denote the input im-
 102 age, output image, and mask, respectively. The prevailing
 103 approach in pixel-space image-to-image frameworks, such
 104 as Palette [19], provides the masked image as a condition-
 105 ing input to the denoiser. However, we find this approach

suboptimal when adapted to the latent space. Specifically, if
 the conditioning input is simply changed to the latent repre-
 sentation of the masked image $\mathcal{E}(\mathbf{m} \odot \mathbf{x})$, the denoiser will
 struggle to preserve structural fidelity around the masked
 region, especially for object removal. This happens even if
 the autoencoder has no difficulty reconstructing the masked
 image, *i.e.*, $\mathcal{D}(\mathcal{E}(\mathbf{m} \odot \mathbf{x})) \approx \mathbf{m} \odot \mathbf{x}$. We attribute this to the
 suboptimal latent representations. Since the autoencoder is
 optimized mainly for reconstruction, its latent space is not
 inherently structured to decouple the original image content
 from the superimposed mask.

To address this, BlazeEdit introduces a secondary, train-
 able encoder $f_\theta(\text{concat}[\mathbf{x}, \mathbf{m}])$ that processes the con-
 catenated input image and mask. The output of f_θ is fed
 to the denoiser as a conditioning input. Unlike the frozen,
 reconstruction-oriented autoencoder, f_θ is jointly optimized
 with the denoiser to extract task-relevant features that the
 model can more easily leverage. We find this joint training
 critical for maintaining fine-grained details and structural
 integrity. Furthermore, for object removal tasks where the
 mask specifies the object to be removed, providing f_θ with
 the full input image \mathbf{x} allows the model to more effectively
 infer and eliminate shadows cast by the object.

Efficient Denoiser and Lightweight Decoder. The de-
 noiser follows the U-ViT architecture [10]. Specifically, we
 employ ResNet [7] blocks at higher resolutions to maintain
 the spatial inductive bias while saving memory, and self-
 attention [23] blocks at lower resolutions to capture long-
 range dependencies and improve accelerator utilization. To
 further reduce model size and inference latency, we prune
 the width and depth of the decoder \mathcal{D} . We then train it for
 image reconstruction with the encoder \mathcal{E} frozen, following
 MobileDiffusion [32]. The resulting lightweight decoder
 has only 6M parameters, and shows minimal degradation in
 reconstruction quality.

141 3.2. Pretraining via Masked Reconstruction

142 A significant bottleneck in developing a generalist image
143 editor is the scarcity of high-quality, task-specific datasets.
144 Existing literature typically follows one of two paths: (1)
145 adapting a pretrained text-to-image model by adding task-
146 specific LoRAs [11, 32], or (2) training an image-to-image
147 model from scratch but limited to a small number of tasks,
148 such as colorization and JPEG restoration [19], where mas-
149 sive paired data can be easily synthesized. The first ap-
150 proach inevitably involves parameters for text understand-
151 ing and text-image alignment, which are unnecessary for
152 purely image-conditioned tasks. On the other hand, the sec-
153 ond approach struggles to scale to more complex editing
154 tasks where the training data are costly to obtain.

155 To achieve data-efficient scaling across diverse editing
156 tasks while maintaining a compact model footprint, we pro-
157 pose to first pretrain an image-to-image base model from
158 scratch on large-scale task-agnostic datasets, and then fine-
159 tune it jointly on all available task-specific datasets.

160 Inspired by the masked image modeling framework [2,
161 8], we use masked reconstruction as our pretraining objec-
162 tive. We find it critical to employ a diverse set of masks,
163 including random patches, geometric shapes, and strokes
164 within the image, and paddings on the boundary of the im-
165 age. This not only encourages the model to learn an expres-
166 sive image representation, but also equips the model with
167 core inpainting and outpainting capabilities, building a gen-
168 eralizable foundation for specialized downstream tasks.

169 We directly leverage the images from existing large-scale
170 text-to-image datasets for our pretraining. The loss function
171 can be written as:

$$172 \mathbb{E}_{\mathbf{x}, \mathbf{m}, \epsilon, t} \|\epsilon_{\theta}(\tilde{\mathbf{x}}_t, f_{\theta}(\text{concat}[\mathbf{m} \odot \mathbf{x}, \mathbf{m}]), t) - \epsilon\|^2, \quad (1)$$

173 where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$ denotes the random noise, and $\tilde{\mathbf{x}}_t$ is
174 the noised version of $\mathcal{E}(\mathbf{x})$ at a randomly sampled diffusion
175 step t . The mask \mathbf{m} is randomly generated on the fly. To prevent
176 degenerate solutions caused by information leak, we mask
177 the image \mathbf{x} when feeding it to the trainable encoder f_{θ} .
178 This is only necessary during pretraining.

179 3.3. Multi-Task Finetuning and Distillation

180 Following pretraining on large-scale, general-purpose im-
181 age data, we transition to supervised finetuning and step
182 distillation on task-specific datasets.

183 **Universal Task Signaling.** We conduct multi-task finetun-
184 ing on all available downstream tasks simultaneously, as
185 this enables knowledge transfer across tasks. To help the
186 model distinguish between different tasks, we introduce a
187 universal task signaling mechanism encoded directly within
188 the mask values. Specifically, each task i is assigned a
189 unique numerical constant τ_i , which is used to scale the bi-
190 nary mask \mathbf{m} . This mask-based conditioning signal allows

the model to switch functional modes dynamically with-
out additional parameters. For an input-mask-output triplet
($\mathbf{x}, \mathbf{m}, \mathbf{y}$) belonging to task i , the finetuning objective is for-
mulated as:

$$195 \mathbb{E}_{\epsilon, t} \|\epsilon_{\theta}(\tilde{\mathbf{y}}_t, f_{\theta}(\text{concat}[\mathbf{x}, \tau_i \cdot \mathbf{m}]), t) - \epsilon\|^2, \quad (2)$$

196 where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$ denotes the random noise, and $\tilde{\mathbf{y}}_t$ is
197 the noised version of $\mathcal{E}(\mathbf{y})$ at a randomly sampled diffusion
198 step t .

199 **Adversarial Distribution Matching Distillation.** After
200 supervised finetuning, we distill the resulting model for 2-
201 step inference, significantly reducing latency. We find that
202 distribution matching distillation [30] combined with adver-
203 sarial training [29] works well in this few-step regime.

204 4. Experiments

205 4.1. Tasks and Datasets

206 We first describe the range of image editing tasks supported
207 by BlazeEdit and the task-specific finetuning datasets.

208 **Object Removal** aims to remove a user-specified object
209 from a scene and synthesize a plausible background in its
210 place, helping to reduce visual clutter. We use a manually
211 curated dataset of $\sim 20\text{K}$ image pairs. Following [24], each
212 pair captures a scene before and after an object is physically
213 removed while minimizing other changes.

214 **Outpainting** extends the boundaries of an image, which is
215 useful for changing its aspect ratio (*e.g.*, portrait to land-
216 scape). We use a high-quality subset of $\sim 5\text{K}$ images from
217 our pretraining datasets, filtered for aesthetic and diversity.

218 **Tone Correction** enhances an image by refining its white
219 balance, saturation, and exposure. We synthesize a dataset
220 of $\sim 3\text{M}$ image pairs by a high-performance teacher model.

221 **Relighting** focuses on portrait enhancement by mitigating
222 unfavorable lighting, such as harsh facial shadows. We use
223 a dataset of $\sim 100\text{K}$ image pairs synthesized by applying
224 diverse shadow augmentations to portrait photography [6].

225 **Sticker Generation** stylizes user-specified subjects within
226 an image (*e.g.*, people, pets, and objects), transforming
227 them into high-quality artistic stickers. We synthesize a di-
228 verse dataset of $\sim 100\text{K}$ image pairs using a high-capacity
229 text-to-image generation model adept at identity-preserving
230 stylization.

231 4.2. Main Results

232 **Qualitative Results.** In Figure 2, we present the input
233 images and generation samples from BlazeEdit on all five
234 tasks. BlazeEdit achieves compelling and versatile image-
235 to-image editing. It exhibits strong structural reasoning
236 and semantic preservation, seamlessly handling object and
237 shadow removal, boundary extrapolation, stylization, tone
238 correction, and lighting adjustment.

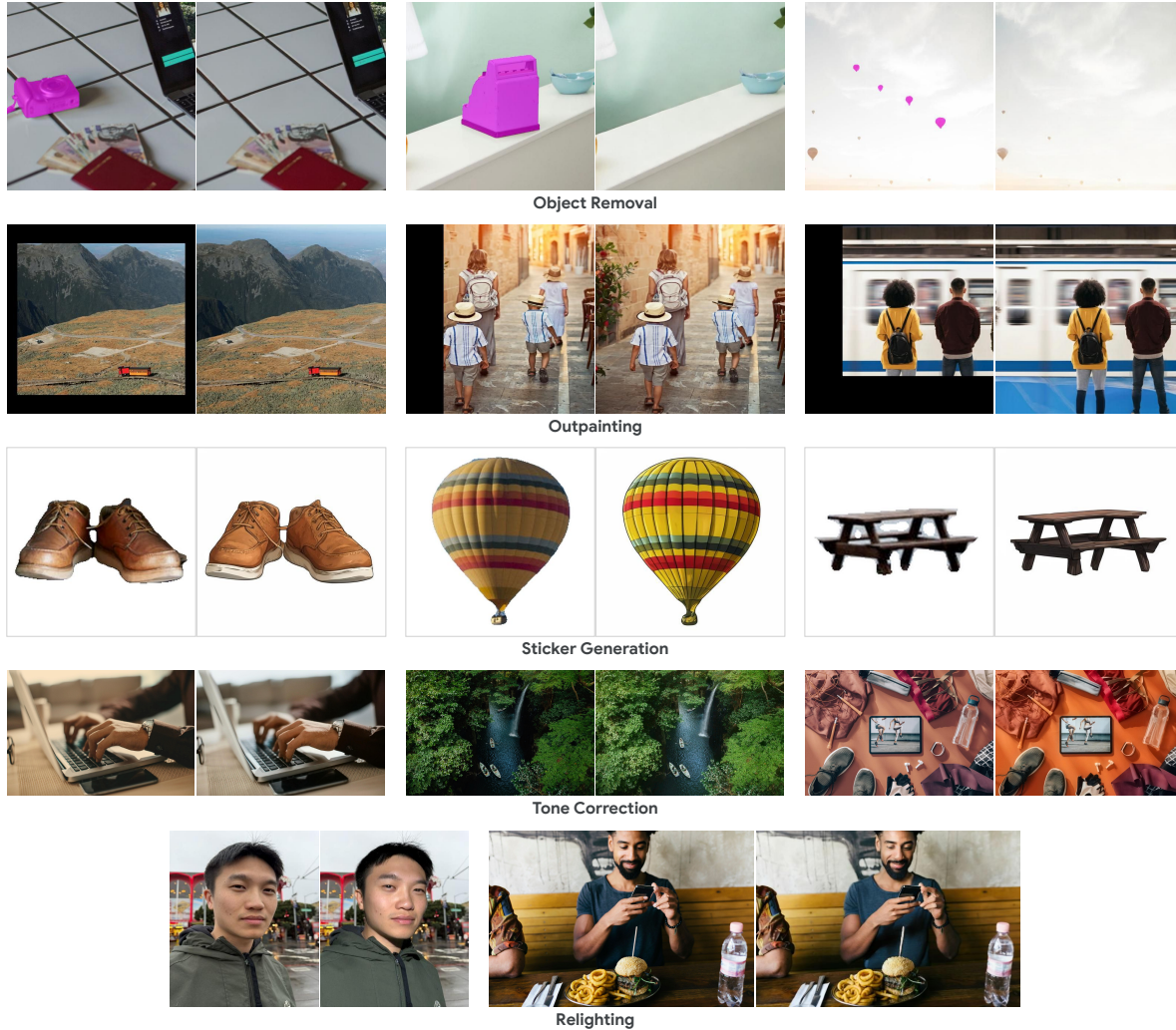


Figure 2. Qualitative results of BlazeEdit on five distinct editing tasks. BlazeEdit achieves compelling and versatile image-to-image editing. It exhibits strong structural reasoning and semantic preservation, seamlessly handling object and shadow removal, boundary extrapolation, stylization, tone correction, and lighting adjustment.

239 **Efficiency Evaluation.** In Table 1, we compare the model
 240 footprint of BlazeEdit with prior works. BlazeEdit achieves
 241 $2\times$ reduction in denoiser parameter count while eliminat-
 242 ing the text encoders, which typically range from 0.1B to
 243 2B parameters. Furthermore, Table 2 provides a break-
 244 down of inference latency on Pixel 10 using Edge TPU.
 245 BlazeEdit completes a full inference pass in just 290ms, en-
 246 abling lightning-fast and privacy-preserving image editing
 247 directly on the edge.

248 5. Conclusion

249 We presented BlazeEdit, a highly efficient image-to-image
 250 diffusion model designed for edge-native generalist image
 251 editing. We successfully consolidated five diverse editing
 252 tasks into a unified 195M-parameter architecture. Our ex-
 253 periments demonstrate that BlazeEdit achieves competitive

Table 1. Model footprint comparison. BlazeEdit achieves $2\times$ re-
 duction in denoiser parameters while eliminating text encoders.

| Model | Text Encoder | Denoiser #Params |
|----------------------|------------------------------------|------------------|
| SnapFusion [13] | CLIP-ViT-H | 848M |
| MobileDiffusion [32] | CLIP-ViT-L | 386M |
| SnapGen [3] | CLIP-ViT-L, CLIP-ViT-G, Gemma-2-2B | 379M |
| BlazeEdit | None | 189M |

Table 2. Inference latency measurements. BlazeEdit completes
 a full inference pass in just 290ms, enabling lightning-fast and
 privacy-preserving image editing directly on the edge.

| Device | Encoder | Decoder | Denoiser (2 steps) | Overall |
|---------------------|---------|---------|--------------------|---------|
| Pixel 10 (Edge TPU) | 45ms | 55ms | 190ms | 290ms |

results while delivering highly interactive user experience
 on mobile hardware.

254
 255

256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312**References**

- [1] Michael Samuel Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. In *International Conference on Learning Representations*, 2023. 1
- [2] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. MaskGIT: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022. 3
- [3] Jierun Chen, Dongting Hu, Xijie Huang, Huseyin Coskun, Arpit Sahni, Aarush Gupta, Anujraaj Goyal, Dishani Lahiri, Rajesh Singh, Yerlan Idelbayev, Junli Cao, Yanyu Li, Kwang-Ting Cheng, S.-H. Gary Chan, Mingming Gong, Sergey Tulyakov, Anil Kag, Yanwu Xu, and Jian Ren. SnapGen: Taming high-resolution text-to-image models for mobile devices with efficient architectures and training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7997–8008, 2025. 1, 2, 4
- [4] 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, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis. In *International Conference on Machine Learning*, 2024. 1
- [5] Tsu-Jui Fu, Yusu Qian, Chen Chen, Wenze Hu, Zhe Gan, and Yinfei Yang. UniVG: A generalist diffusion model for unified image generation and editing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17160–17170, 2025. 2
- [6] David Futschik, Kelvin Ritland, James Vecore, Sean Fanello, Sergio Orts-Escolano, Brian Curless, Daniel Šykora, and Rohit Pandey. Controllable light diffusion for portraits. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8412–8421, 2023. 3
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 2
- [8] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 3
- [9] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, 2020. 1
- [10] Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. simple diffusion: End-to-end diffusion for high resolution images. In *International Conference on Machine Learning*, 2023. 2
- [11] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. 3
- [12] Duong H. Le, Tuan Pham, Sangho Lee, Christopher Clark, Aniruddha Kembhavi, Stephan Mandt, Ranjay Krishna, and Jiasen Lu. One diffusion to generate them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2671–2682, 2025. 2
- [13] Yanyu Li, Huan Wang, Qing Jin, Ju Hu, Pavlo Chemerys, Yun Fu, Yanzhi Wang, Sergey Tulyakov, and Jian Ren. SnapFusion: Text-to-image diffusion model on mobile devices within two seconds. In *Advances in Neural Information Processing Systems*, 2023. 1, 2, 4
- [14] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *International Conference on Learning Representations*, 2023. 1
- [15] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *International Conference on Learning Representations*, 2023. 1
- [16] Can Qin, Shu Zhang, Ning Yu, Yihao Feng, Xinyi Yang, Yingbo Zhou, Huan Wang, Juan Carlos Niebles, Caiming Xiong, Silvio Savarese, Stefano Ermon, Yun Fu, and Ran Xu. UniControl: A unified diffusion model for controllable visual generation in the wild. In *Advances in Neural Information Processing Systems*, 2023. 2
- [17] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, 2021. 1
- [18] 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*, pages 10684–10695, 2022. 2
- [19] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. 2, 3
- [20] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, 2015. 1
- [21] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021. 1
- [22] Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison,

- 370 Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, 428
 371 Andy Coenen, Anthony Laforge, Antonia Paterson, Ben 429
 372 Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, 430
 373 Chintu Kumar, Chris Perry, Chris Welty, Christopher A. 431
 374 Choquette-Choo, Danila Sinopalnikov, David Weinberger, 432
 375 Dimple Vijaykumar, Dominika Rogozińska, Dustin Her- 433
 376 bison, Elisa Bandy, Emma Wang, Eric Noland, Erica 434
 377 Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, 435
 378 Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Mart- 436
 379 ins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen 437
 380 Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack 438
 381 Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng 439
 382 Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, 440
 383 Joe Fernandez, Joost van Amersfoort, Josh Gordon, Josh 441
 384 Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed, 442
 385 Kartikeya Badola, Kat Black, Katie Millican, Keelin Mc- 443
 386 Donell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, 444
 387 Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena 445
 388 Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini 446
 389 Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, 447
 390 Machel Reid, Manvinder Singh, Mark Iverson, Martin 448
 391 Görner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt 449
 392 Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, 450
 393 Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk 451
 394 Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, 452
 395 Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta 453
 396 Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, 454
 397 Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, 455
 398 Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona 456
 399 Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir 457
 400 Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, 458
 401 Sara Mc Carthy, Sarah Cogan, Sarah Perrin. Sébastien M. R. 459
 402 Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, 460
 403 Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, 461
 404 Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, 462
 405 Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh Meshram, 463
 406 Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming 464
 407 Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, 465
 408 Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, 466
 409 Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli 467
 410 Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, 468
 411 D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol 469
 412 Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, 470
 413 Clement Farabet, Elena Buchatskaya, Sebastian Borgeaud, 471
 414 Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert 472
 415 Dadashi, and Alek Andreev. Gemma 2: Improving open 473
 416 language models at a practical size, 2024. 1
- [23] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszko- 474
 417 reit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia 475
 418 Polosukhin. Attention is all you need. In *Advances in neural 476
 419 information processing systems*, 2017. 2
- [24] Daniel Winter, Matan Cohen, Shlomi Fruchter, Yael Pritch, 477
 421 Alex Rav-Acha, and Yedid Hoshen. ObjectDrop: Bootstrap- 478
 422 ping counterfactuals for photorealistic object removal and 479
 423 insertion. In *European Conference on Computer Vision*, pages 480
 424 112–129. Springer, 2024. 3
- [25] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan 481
 426 Gao, Kun Yan, Sheng ming Yin, Shuai Bai, Xiao Xu, Yilei 482
 427 Chen, Yuxiang Chen, Zecheng Tang, Zekai Zhang, Zhengyi 483
 428 Wang, An Yang, Bowen Yu, Chen Cheng, Dayiheng Liu, De- 484
 429 qing Li, Hang Zhang, Hao Meng, Hu Wei, Jingyuan Ni, Kai 485
 430 Chen, Kuan Cao, Liang Peng, Lin Qu, Minggang Wu, Peng 486
 431 Wang, Shuting Yu, Tingkun Wen, Wensen Feng, Xiaoxiao 487
 432 Xu, Yi Wang, Yichang Zhang, Yongqiang Zhu, Yujia Wu, 488
 433 Yuxuan Cai, and Zenan Liu. Qwen-Image technical report, 489
 434 2025. 1, 2
- [26] Bin Xia, Yuechen Zhang, Jingyao Li, Chengyao Wang, 490
 435 Yitong Wang, Xinglong Wu, Bei Yu, and Jiaya Jia. 491
 436 DreamOmni: Unified image generation and editing. In *Pro- 492
 437 ceedings of the IEEE/CVF Conference on Computer Vision 493
 438 and Pattern Recognition*, pages 28533–28543, 2025. 494
- [27] Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xin- 495
 439 grun Xing, Ruiran Yan, Chaofan Li, Shuting Wang, Tiejun 496
 440 Huang, and Zheng Liu. OmniGen: Unified image genera- 497
 441 tion. In *Proceedings of the IEEE/CVF Conference on Com- 498
 442 puter Vision and Pattern Recognition*, pages 13294–13304, 499
 443 2025. 444
- [28] Xingqian Xu, Zhangyang Wang, Gong Zhang, Kai Wang, 500
 445 and Humphrey Shi. Versatile Diffusion: Text, images and 501
 446 variations all in one diffusion model. In *Proceedings of the 502
 447 IEEE/CVF International Conference on Computer Vision*, 503
 448 pages 7754–7765, 2023. 2
- [29] Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, 504
 449 Eli Shechtman, Fredo Durand, and Bill Freeman. Improved 505
 450 distribution matching distillation for fast image synthesis. In 506
 451 *Advances in neural information processing systems*, 2024. 3
- [30] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shecht- 507
 452 man, Fredo Durand, William T Freeman, and Taesung Park. 508
 453 One-step diffusion with distribution matching distillation. In 509
 454 *Proceedings of the IEEE/CVF Conference on Computer Vi- 510
 455 sion and Pattern Recognition*, pages 6613–6623, 2024. 3
- [31] Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin 511
 456 Bao, Shaozhe Hao, Lu Yuan, and Kwan-Yee K Wong. Uni- 512
 457 ControlNet: All-in-one control to text-to-image diffusion 513
 458 models. In *Advances in neural information processing sys- 514
 459 tems*, 2023. 2
- [32] Yang Zhao, Yanwu Xu, Zhisheng Xiao, Haolin Jia, and 515
 460 Tingbo Hou. MobileDiffusion: Instant text-to-image genera- 516
 461 tion on mobile devices. In *European Conference on Com- 517
 462 puter Vision*, pages 225–242. Springer, 2024. 1, 2, 3, 4