
Appendix to “FreCaS: Efficient Higher-Resolution Image Generation via Frequency-aware Cascaded Sampling”

In this appendix, we provide the following materials:

- A Details of timestep shifting in the transition process (referring to Sec. 3.2 in the main paper);
- B The detailed settings of FreCaS on $\times 4$ and $\times 16$ generation for SD2.1, SDXL and SD3 (referring to Sec. 4.1 in the main paper);
- C [Results of user studies and non-reference image quality assessment \(NR-IQA\)](#) (referring to Sec. 4.1 in the main paper);
- D [Comparison with training-based methods and super-resolution methods](#) (referring to Sec. 4.2 in the main paper);
- E [More visual results and visual comparisons](#) (referring to Sec. 4.2 in the main paper);
- F Experimental results of generation of SD3 (referring to Sec. 4.3 in the main paper);
- G [Ablation studies on individual components of FreCaS and inference schedule](#) (referring to Sec. 4.4 in the main paper).

A SHIFTING TIMESTEP IN THE TRANSITION PROCESS

As mentioned in Sec. 3.2 of the main paper, FreCaS employs a five-step transition process to transform the last latent in the current stage $z_L^{s_{i-1}}$ to the first latent in the next stage $z_F^{s_i}$. In addition to changing the resolution, we adjust the timestep from L to F to ensure that the signal-to-noise ratio (SNR) (Kingma et al., 2021) could be a constant in the transition process. Given a state $z_t = \sqrt{\alpha_t}z_0 + \sqrt{1 - \alpha_t}\epsilon$ at timestep t , the SNR is defined as $\text{SNR}(z_t) = \frac{\alpha_t}{1 - \alpha_t}$, where $\alpha_1, \dots, \alpha_T$ represent the noise schedule, and ϵ is Gaussian noise. It has been found (Hoogeboom et al., 2023; Chen, 2023) that the SNR maintains a consistent ratio across resolutions for diffusion models using the same noise schedule:

$$\text{SNR}(z_t^s) = \text{SNR}(z_t^{\hat{s}}) \cdot \left(\frac{s}{\hat{s}}\right)^\gamma,$$

where s and \hat{s} denote different resolutions. The value of γ is typically set to 2.

Teng et al. (2024) and Gu et al. (2023) proposed to redesign the noise schedule to keep SNR consistent when changing the resolutions of intermediate states. Since the pre-trained diffusion models have fixed noise schedules, in this paper we adjust the timestep, instead of the noise schedule, to ensure consistent SNR between $z_L^{s_{i-1}}$ and $z_F^{s_i}$:

$$\text{SNR}(z_L^{s_{i-1}}) = \text{SNR}(z_F^{s_i}) \Rightarrow F = \alpha^{-1} \left(\frac{\left(\frac{s_{i-1}}{s_i}\right)^\gamma \cdot \alpha_L}{1 + \left(\left(\frac{s_{i-1}}{s_i}\right)^\gamma - 1\right) \cdot \alpha_L} \right), \quad (1)$$

where α^{-1} is the inverse function of α_t . Proper adjustment of γ can yield additional improvements.

Besides, SD3 (Esser et al., 2024) employs a similar formula to shift the timestep when varying resolutions:

$$F = \frac{\sqrt{\frac{s_i}{s_{i-1}}} \cdot L}{1 + \left(\sqrt{\frac{s_i}{s_{i-1}}} - 1\right) \cdot L}. \quad (2)$$

B EXPERIMENTAL SETTING DETAILS

The experimental setting details of our FreCaS are listed in Table 1.

C RESULTS OF USER STUDIES AND NR-IQA METRICS

We have (a) conducted user studies and (b) employed non-reference image quality assessment (NR-IQA) metrics to further assess the performance of FreCaS and its competing methods.

Table 1: Detailed settings of FreCaS on the experiments. N denotes the count of additional stages. “Steps” presents the sampling steps in each stage. L presents the timestep of last latent in each stage except for the final one. γ denotes the SNR ratio in the transition process. w_l , w_h and w_c are the hyper-parameters of the proposed FA-CFG and CA-maps re-utilization.

		$N + 1$	Steps	L	γ	w_l	w_h	w_c
SD2.1	$\times 4$	2	40,10	100	3.0	7.5	45.0	0.6
	$\times 16$	3	30,10,10	200,200	3.0	7.5	35.0	0.4
SDXL	$\times 4$	2	40,10	200	1.5	7.5	35.0	0.6
	$\times 16$	3	30,5,15	400,200	2.0	7.5	35.0	0.6
SD3	$\times 4$	2	20,8	50	-	7.0	35.0	0.5

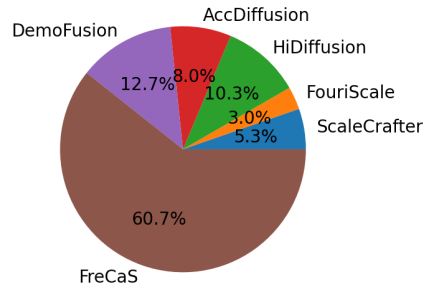


Figure 1: User study results on $\times 4$ generation of SDXL.

Table 2: NR-IQA metrics on $\times 4$ and $\times 16$ generation of SDXL.

Methods	$\times 4$			$\times 16$		
	clipiqa \uparrow	niqe \downarrow	musiq \uparrow	clipiqa \uparrow	niqe \downarrow	musiq \uparrow
DirectInference	0.522	4.167	53.98	0.469	4.370	29.00
AttnEntropy	0.547	4.210	54.87	0.528	4.614	27.98
ScaleCrafter	0.664	3.577	61.12	0.618	3.783	36.00
FouriScale	0.662	3.580	60.77	0.612	3.791	35.52
HiDiffusion	0.690	4.049	61.69	0.574	7.348	36.71
AccDiffusion	0.627	3.641	57.02	0.626	3.587	31.83
DemoFusion	0.651	3.410	58.98	0.637	3.376	33.46
Ours	0.668	3.391	63.10	0.646	3.367	37.33

C.1 USER STUDIES

For the user studies, we compare FreCaS with ScaleCrafter, FouriScale, HiDiffusion, DemoFusion, and AccDiffusion on 2048×2048 image generation using SDXL. We randomly selected 30 prompts and generated one image per method for each prompt, creating 30 sets of images. Ten volunteers participated in the test, and they were asked to select the image with the best details and reasonable semantic layout from each set. The results are shown in Figure 1. We can see that FreCaS significantly outperforms other methods, with 60% votes as the best method. DemoFusion, AccDiffusion, and HiDiffusion perform similarly, with each having about 10% of the votes. In contrast, FouriScale and ScaleCrafter have the fewest votes, about 5% each.

C.2 NR-IQA METRICS

For the NR-IQA metrics, we employ CLIPIQA (Wang et al., 2023), NIQE (Mittal et al., 2012), and MUSIQ (Ke et al., 2021) on $\times 4$ and $\times 16$ image generations with SDXL. The results are presented in

108
 109
 110
 111
 112
 113
 114
 115
 116
 117
 118
 119
 120
 121
 122
 123
 124
 125
 126
 127
 128
 129
 130
 131
 132
 133
 134
 135
 136
 137
 138
 139
 140
 141
 142
 143
 144
 145
 146
 147
 148
 149
 150
 151
 152
 153
 154
 155
 156
 157
 158
 159
 160
 161

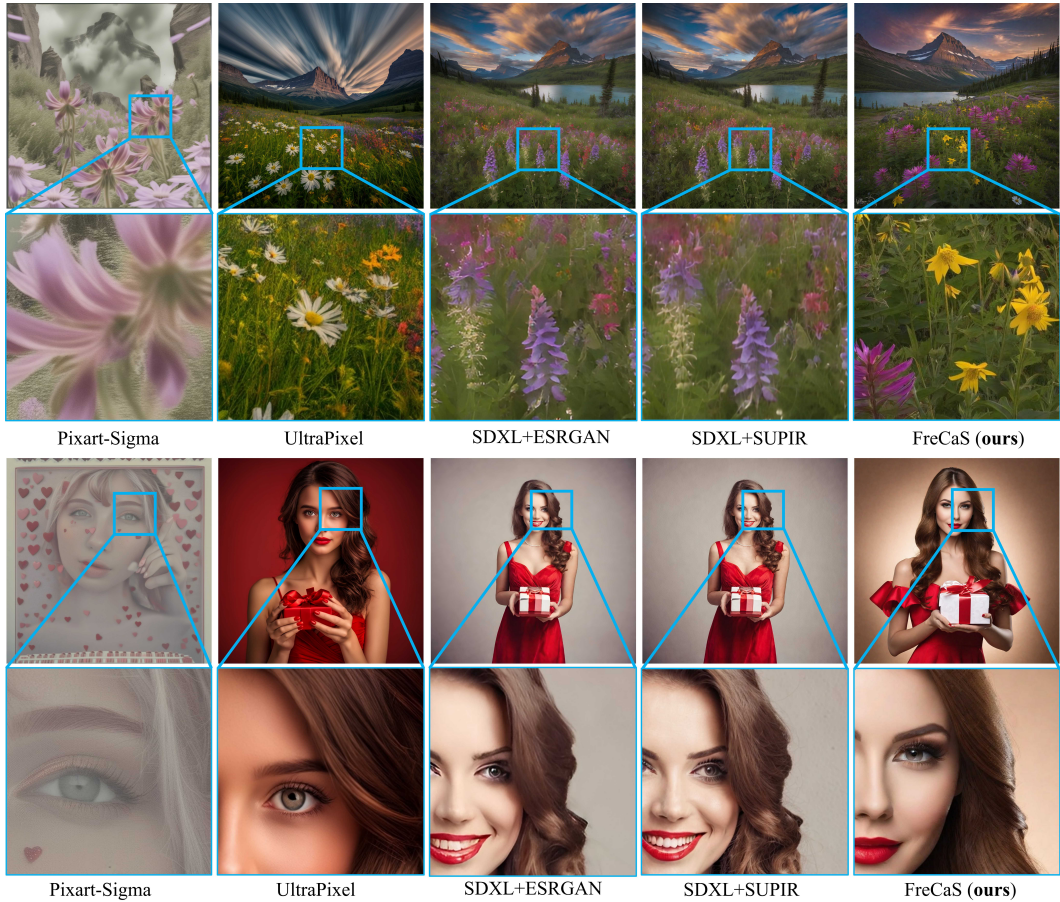


Figure 2: Visual comparison with training-based methods and super-resolution methods on $\times 4$ generation of SDXL.

Table 2. Our FreCaS consistently outperforms all the other methods. For example, on $\times 4$ generation, FreCaS achieves a CLIPQA score of 0.668, a NIQE score of 3.391, and a MUSIQ score of 63.10, compared to 0.651, 3.410, and 58.98 for DemoFusion. On $\times 16$ generation, FreCaS achieved a CLIPQA score of 0.646, a NIQE score of 3.367, and a MUSIQ score of 37.33, compared to 0.626, 3.587, and 31.83 for AccuDiffusion. Notably, FreCaS only lags behind HiDiffusion on the CLIPQA metric in $\times 4$ image generation.

D COMPARISON WITH TRAINING-BASED METHODS AND SUPER-RESOLUTION METHODS

We conducted additional experiments comparing FreCaS with training-based methods (Pixart-Sigma (Chen et al., 2024) and UltraPixel (Ren et al., 2024)) and super-resolution methods (ESRGAN (Wang et al., 2021) and SUPIR (Yu et al., 2024)). To ensure fair comparisons, we set the model precision to fp16 (bf16 for UltraPixel, as recommended by the authors) and use the DDIM sampler for diffusion-based methods. For Pixart-Sigma, we can only report its results for 2048×2048 image generation since its 4K model is not available. The quantitative results are summarized in Table 3.

From Table 3, we can see that FreCaS outperforms Pixart-Sigma and UltraPixel in most metrics. For example, FreCaS achieves an FID score of 16.48 and an IS score of 17.18, compared to 26.1 and 14.44 of Pixart-Sigma, and 25.56 and 17.11 of UltraPixel on the $\times 4$ image generation task. This is because Pixart-Sigma, as acknowledged by the authors, heavily relies on the advanced samplers (see <https://github.com/PixArt-alpha/PixArt-sigma/issues/65>) so that the results are not very stable.

Table 3: Comparison with training-based methods and super-resolution methods on $\times 4$ and $\times 16$ generation of SDXL.

	Methods	FID \downarrow	FID $_p\downarrow$	IS \uparrow	IS $_p\uparrow$	CLIP SCORE \uparrow	Latency(s) \downarrow
$\times 4$	Pixart-Sigma	26.11	38.58	14.44	14.45	28.10	71.45
	UltraPixel	25.56	19.95	17.11	17.10	33.17	41.70
	SDXL+ESRGAN	13.03	18.10	17.30	16.58	34.13	6.36
	SDXL+SUPIR	12.08	17.31	17.57	17.12	34.16	105.5
	Ours	16.48	17.91	17.18	17.31	33.28	13.84
$\times 16$	UltraPixel	51.43	45.88	12.48	13.73	33.07	162.4
	SDXL+ESRGAN	45.86	43.10	12.94	13.48	33.44	7.25
	SDXL+SUPIR	43.94	39.35	13.22	14.37	33.49	512.4
	Ours	42.75	39.82	12.68	14.16	33.03	85.87

UltraPixel, while achieving comparable performance to DemoFusion, still lags behind FreCaS in most metrics. Besides, the two methods are much slower than our FreCaS.

For SR-based methods, FreCaS may have lower FID, IS, and CLIP scores than SDXL+ESRGAN. This is because SR methods are designed to strictly adhere to low-resolution inputs, while these metrics (FID, IS, and CLIP) evaluate images by downsampling them to low resolution, which cannot well reflect the quality of generated high-resolution images. However, FreCaS significantly outperforms SDXL+ESRGAN in FID $_p$ and IS $_p$. Specifically, FreCaS achieves an FID $_p$ score of 39.82 and an IS $_p$ score of 14.16, compared to 43.10 and 13.48 of SDXL+ESRGAN on $\times 16$ image generation. This indicates its superior ability to generate high-resolution local details. This observation is consistent with the findings in the DemoFusion paper. Additionally, SDXL+SUPIR outperforms FreCaS in FID $_p$ and IS $_p$, but at the cost of much longer inference latency (85.87 seconds for FreCaS vs. 512.4 seconds for SDXL+SUPIR on $\times 16$ image generations).

We have provided some visual comparisons in Figure 2. One can see that FreCaS demonstrates better visual quality than either training-based or SR-based methods in high-resolution image generation, such as the more vivid and clearer flowers, hairs and the more natural color of lips.

E MORE VISUAL RESULTS

E.1 MORE VISUAL RESULTS

Figure 3 illustrates more visual results of FreCaS, including those with varying aspect ratios. From top to bottom, and left to right, the prompts used in examples are: 1. “Beautiful winter wallpapers.” 2. “A regal queen adorned with jewels.” 3. “A majestic phoenix, wings ablaze, rising from ashes, the flames casting a warm glow.” 4. “Lady in Red oil portrait painting won the John Singer Sargent People’s award.” 5. “Star of the day – Actress Evelyn Laye - 1917.” 6. “Photograph - Clouds Over Daicey Pond by Rick Berk.” 7. “little-boy-with-large-bulldog-in-a-garden-france.” 8. “03-Brussels-Maja-Wronska-Travels-Architecture-Paintings.” 9. “Red Fox Pup Print by William H. Mullins.” 10. “Lovely Illustrations Of Cityscapes Inspired By Southeast Asia Malaysian digital illustrator Chong Fei Giap’s illustrations of cityscapes are lovely and inspiring. Fantasy Landscape, Landscape Art, Illustrator, Japon Tokyo, Animation Background, Art Background, Matte Painting, Anime Scenery, Jolie Photo.” 11. “A plate with creamy chicken and vegetables, a side of onion rings, a cup of coffee and a slice of cheesecake.” 12. “Hyper-Realistic Portrait of Redhead Girl Drawn with Bic Pens.”

To further validate the performance of FreCaS in real-world application scenarios, we have provided additional visual results in three categories:

- **Simple scenes.** These images typically contain a single object in a realistic style. We display images of people, animals, landscapes, buildings, and other common objects. The visual results for this group are presented in Figure 4.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269



Figure 3: Visual results of FreCaS on SDXL. Please zoom-in for better view.

- **Various styles.** This group showcases images in different artistic styles, including oil painting, pencil sketch, ink wash, watercolor, and poster art. The results are shown in the first two rows of Figure 5.

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
313
314
315
316
317
318
319
320
321
322
323



Figure 4: More visual results on simple scenes.

- **Complex scenes.** These images contain multiple objects or have intricate textures. The results are displayed in the bottom two rows of are presented in Figure 5.

From these visual results, it is evident that FreCaS consistently generates high-quality images across various styles and contents, demonstrating the capability of FreCaS in real-world applications.

z

E.2 MORE VISUAL COMPARISONS

We show more visual comparisons in Figure 6. From top to bottom, the prompts used in the four groups of examples are: 1. “A small den with a couch near the window.” 2. “A painting of a candlestick holder with a candle, several pieces of fruit and a vase, with a gold frame around the painting.” 3. “A noble knight, riding a white horse, the castle gates opening.” 4. “Mystical Landscape Digital Art - Lonely Tree Idyllic Winterlandscape by Melanie Viola.”

We have provided more 4K visual comparisons under realistic scenes in Figure 7. As can be seen, our FreCaS consistently delivers better results in both image layout and semantic details.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377



Figure 5: More visual results of various styles (top two rows) and complex scenes (bottom two rows).



425 Figure 6: Visual comparisons on $\times 4$ and $\times 16$ experiments of SD2.1 and SDXL. Please zoom-in for
426 better view.

427
428
429 **F EXPERIMENTS ON SD3**

430 In this section, we present the results of the $\times 4$ generation experiments on SD3. SD3 employs
431 a transformer-based denoising network. It eliminates all convolutional layers, thereby preventing

432
 433
 434
 435
 436
 437
 438
 439
 440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451
 452
 453
 454
 455
 456
 457
 458
 459
 460
 461
 462
 463
 464
 465
 466
 467
 468
 469
 470
 471
 472
 473
 474
 475
 476
 477
 478
 479
 480
 481
 482
 483
 484
 485

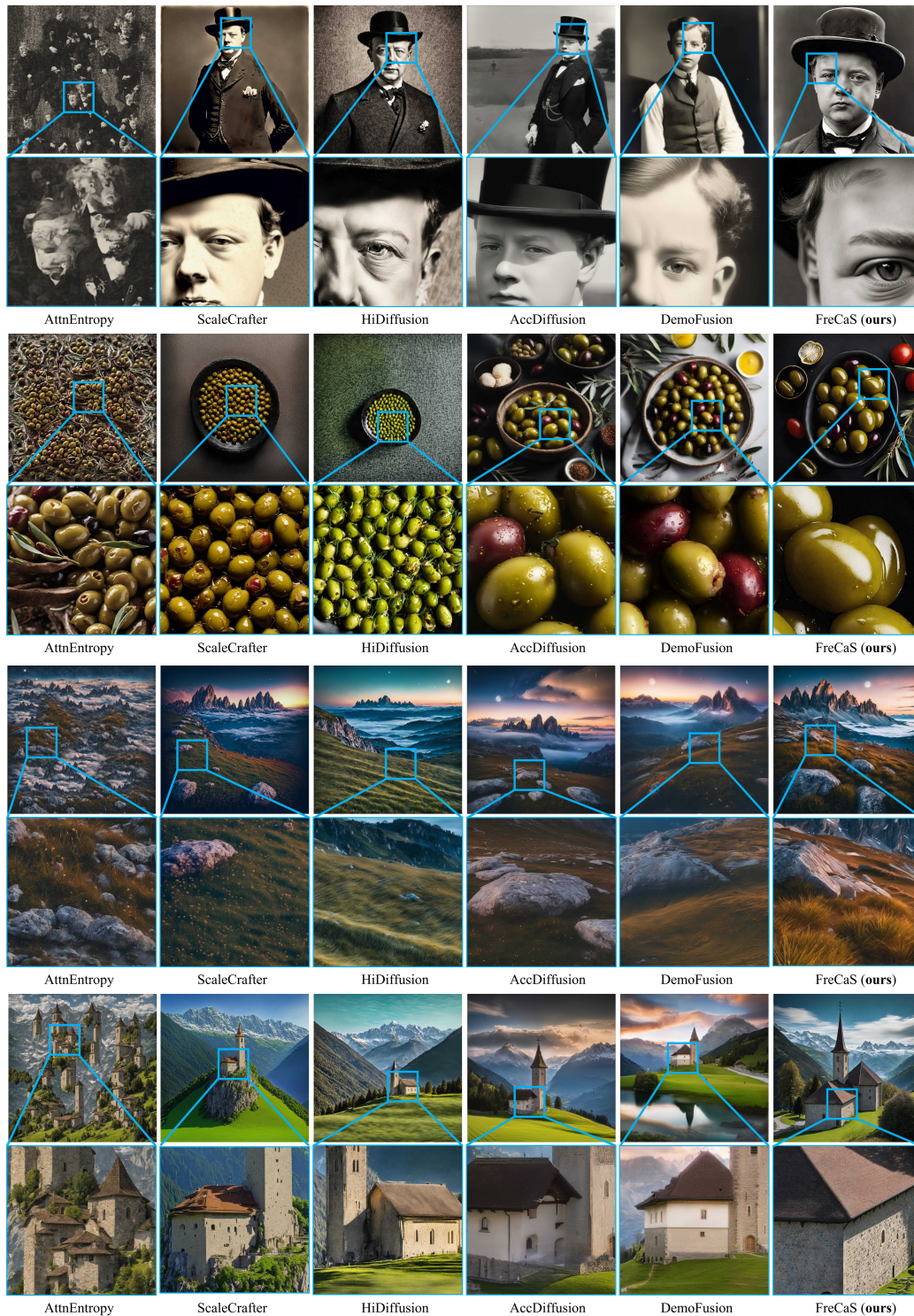


Figure 7: More 4K comparisons in realistic styles. From top to bottom, the prompts are “Young winston churchill.”, “Olive food photography.”, “Mountains in fog at beautiful night. Dreamy landscape with mountain peaks, stones, grass, blue sky with blurred low clouds, stars and moon. Rocks at dusk.” and “Image Church Switzerland towers San Romerio Nature Mountains Scenery Made of stone Tower mountain landscape photography.”

486
487
488
489
490
491
492
493
494

Table 4: Experiments on $\times 4$ generation of SD3.

Methods	FID _b ↓	FID _p ↓	IS↑	IS _p ↑	CLIP SCORE↑	Latency (s)↓	SpeedUP↑
DirectInference	35.68	45.35	12.52	12.60	31.45	38.53	1×
Demodiffusion	15.19	44.34	17.84	14.99	31.09	63.33	0.61×
Ours	9.76	26.62	17.83	16.72	31.17	15.94	2.42×

495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518

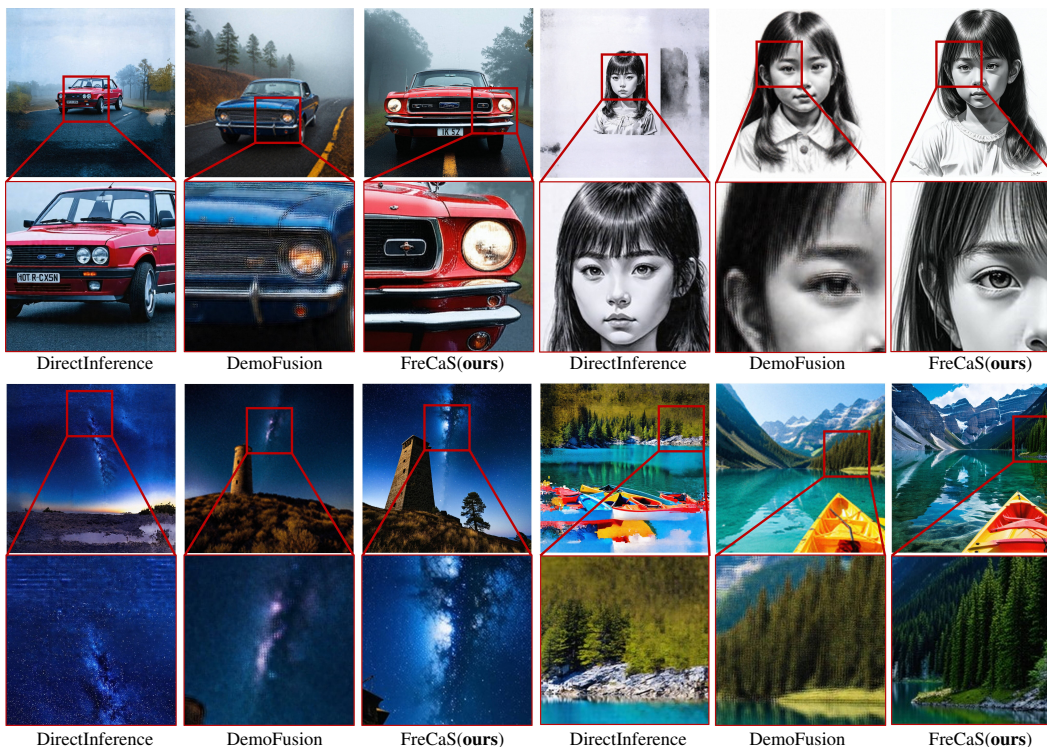


Figure 8: Visual comparison on $\times 4$ experiments of SD3. From top to bottom, from left to right, the prompts used in the four groups of examples are: 1. “Car Photograph - Ford In The Fog by Debra and Dave Vanderlaan.” 2. “Rupert Young is Sir Leon in Merlin season 5 copy.” 3. “Watchtower, Shooting Star & Milky Way, Gualala, CA.” 4. “Colorful Autumn in Mount Fuji, Japan - Lake Kawaguchiko is one of the best places in Japan to enjoy Mount Fuji scenery of maple leaves changing color giving image of those leaves framing Mount Fuji.”. Zoom-in for better view.

525
526
527
528
529
530
531

the application of many existing methods, such as ScaleCrafter and FouriScale. Besides, SD3 exhibits fine details in the central region but shows corrupted textures in the surrounding regions (see Figure 8). This issue with the image layout also significantly impacts the performance of other methods, such as DemoFusion. Therefore, we only compare our FreCaS with DirectInference and DemoFusion. Table 4 and Figure 8 present the quantitative and qualitative results, respectively.

532
533
534
535
536
537
538
539

From Table 4, it is evident that FreCaS achieves superior performance in terms of image quality and inference speed. Specifically, FreCaS achieves the best results on FID_b, FID_p, IS, and IS_p, and only slightly lags behind DirectInference in terms of CLIP score. Moreover, FreCaS generates a 2048×2048 image in about 16 seconds, achieving a speed-up of $2.42\times$ and $3.97\times$ compared to DirectInference and DemoFusion, respectively. Figure 8 illustrates the generated images. Directly employing the pre-trained SD3 model to generate higher-resolution images, DirectInference leads to unreasonable image layout with the surrounding parts being corrupted, such as the road and trees. The results of DemoFusion exhibits strange artifacts, such as the car faces and eyes. In contrast, our FreCaS successfully maintains the natural image structure while obtaining fine details.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

Table 5: Ablation studies on 2048×2048 generation of SDXL.

Model	cascaded framework	FA-CFG	CA-reuse	FID \downarrow	FID $_p \downarrow$	IS \uparrow	IS $_p \uparrow$	CLIP SCORE \uparrow	Latency (s)
#1				39.14	29.71	11.52	14.60	32.51	34.10
#2	✓			17.62	20.49	17.01	16.54	33.24	13.71
#3	✓	✓		16.62	17.91	17.16	16.82	33.34	13.74
#4	✓	✓	✓	16.48	17.91	17.18	17.31	33.28	13.84

Table 6: Ablation studies on N in FreCaS.

N	resolutions	FID $_b \downarrow$	FID $_p \downarrow$
0	2048	43.83	29.71
1	1024 \rightarrow 2048	12.63	17.91
2	1024 \rightarrow 1536 \rightarrow 2048	41.36	28.68

Table 7: Ablation studies on L in FreCaS.

L	FID $_b \downarrow$	FID $_p \downarrow$
0	12.57	18.20
100	12.69	18.10
200	12.63	17.91
300	13.30	18.57
400	13.34	18.62

G ABLATION STUDIES ON INDIVIDUAL COMPONENTS AND INFERENCE SCHEDULE

We further conduct ablation studies to verify the effectiveness of each components and the settings of inference schedule of our FreCaS.

G.1 EFFECTIVENESS OF EACH COMPONENT

To better verify the effectiveness of each component of FreCaS, we conducted more ablation studies on our proposed cascaded framework, FA-CFG, and CA-reuse strategies. The results are shown in Table 5. One can see that our cascaded framework significantly outperforms the baseline, with a decrease of 22.52 in the FID score and a reduction of 20.39 seconds in latency. This demonstrates the high efficiency of our proposed cascaded framework. Our FA-CFG strategy improves both FID and IS scores and shows substantial improvement in FID $_p$, demonstrating its effectiveness in generating realistic image details. The CA-reuse strategy further enhances IS $_p$, indicating its effectiveness in improving semantic appearance. Moreover, these strategies introduce minimal additional latency.

G.2 EXPERIMENTS ON INFERENCE SCHEDULE

In this section, we conduct experiments on the selection of N (number of additional stages) and L (the timestep of last latent in each stage). The two factors are employed to adjust the inference schedule of our FreCaS. We reports the scores of FID $_b$ and FID $_p$ by varying the two factors in Table 6 and Table 7, respectively.

Choice of N . From Table 6, we see that $N = 1$ achieves an FID $_b$ score of 12.63 and an FID $_p$ score of 17.91, significantly better than $N = 0$ and $N = 2$ in the $\times 4$ generation task for SDXL. This could be attributed to the fact that a larger value of N introduces more transition steps, which can lead to much information loss. Conversely, a smaller value of N reduces the effectiveness of FreCaS, degenerating it to the DirectInference method.

594 **Choice of L .** From Table 7, we can see that a smaller L improves FID_b score but deteriorates FID_p .
595 This is because the details generated at lower resolutions conflict with those at higher resolutions.
596 Thus, we set L to 200 to avoid generating excessive unwanted details in the early stages.
597

598 REFERENCES

- 599
600 Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping
601 Luo, Huchuan Lu, and Zhenguo Li. Pixart- σ : Weak-to-strong training of diffusion trans-
602 former for 4k text-to-image generation. *arXiv preprint arXiv:2403.04692*, 2024.
603
604 Ting Chen. On the importance of noise scheduling for diffusion models. *arXiv preprint*
605 *arXiv:2301.10972*, 2023.
606
607 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam
608 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for
609 high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*,
610 2024.
611
612 Jiatao Gu, Shuangfei Zhai, Yizhe Zhang, Miguel Ángel Bautista, and Joshua M Susskind. f-dm: A
613 multi-stage diffusion model via progressive signal transformation. In *The Eleventh International*
614 *Conference on Learning Representations*, 2023.
615
616 Emiel Hoogetboom, Jonathan Heek, and Tim Salimans. simple diffusion: End-to-end diffusion for
617 high resolution images. In *International Conference on Machine Learning*, pp. 13213–13232.
618 PMLR, 2023.
619
620 Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image
621 quality transformer. In *Proceedings of the IEEE/CVF International Conference on Computer*
622 *Vision*, pp. 5148–5157, 2021.
623
624 Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Ad-*
625 *vances in neural information processing systems*, 34:21696–21707, 2021.
626
627 Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. Making a “completely blind” image quality
628 analyzer. *IEEE Signal processing letters*, 20(3):209–212, 2012.
629
630 Jingjing Ren, Wenbo Li, Haoyu Chen, Renjing Pei, Bin Shao, Yong Guo, Long Peng, Fenglong
631 Song, and Lei Zhu. Ultrapixel: Advancing ultra-high-resolution image synthesis to new peaks.
632 *arXiv preprint arXiv:2407.02158*, 2024.
633
634 Jiayan Teng, Wendi Zheng, Ming Ding, Wenyi Hong, Jianqiao Wangni, Zhuoyi Yang, and Jie Tang.
635 Relay diffusion: Unifying diffusion process across resolutions for image synthesis. 2024.
636
637 Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. Exploring clip for assessing the look and
638 feel of images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp.
639 2555–2563, 2023.
640
641 Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind
642 super-resolution with pure synthetic data. In *Proceedings of the IEEE/CVF international confer-*
643 *ence on computer vision*, pp. 1905–1914, 2021.
644
645 Fanghua Yu, Jinjin Gu, Zheyuan Li, Jinfan Hu, Xiangtao Kong, Xintao Wang, Jingwen He, Yu Qiao,
646 and Chao Dong. Scaling up to excellence: Practicing model scaling for photo-realistic image
647 restoration in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
Pattern Recognition, pp. 25669–25680, 2024.